

# A novel method of gene regulatory network structure inference from gene knock-out expression data

Xiang Chen <sup>1</sup>, Min Li <sup>1\*</sup>, Ruiqing Zheng <sup>1</sup>, Siyu Zhao <sup>1</sup>, Fang-Xiang Wu<sup>1,2</sup>, Yaohang Li<sup>3</sup> and Jianxin Wang <sup>1</sup>

<sup>1</sup> School of Information Science and Engineering, Central South University, Changsha 410083, China; chenxofhit@gmail.com (X.C.); rqzheng@csu.edu.cn (R.Z.); syzhao@csu.edu.cn (S.Z.); faw341@mail.usask.ca (F.W.); jxwang@mail.csu.edu.cn (J.W.)

<sup>2</sup> Department of Mechanical Engineering and Division of Biomedical Engineering, University of Saskatchewan, Saskatoon, Canada

<sup>3</sup> Department of Computer Science, Old Dominion University, Norfolk, VA, USA; yaohang@cs.odu.edu (Y.L.)

\* Correspondence: limin@mail.csu.edu.cn (M.L.); Tel.: +86-0731-888-30212

Academic Editor: Quan Zou; Le Zhang

Version December 11, 2017 submitted to Int. J. Mol. Sci.

**Abstract:** Inferring gene regulatory networks (GRNs) structure from gene expression data has been a pretty challenging problem in systems biology. It is critical to identify complicated regulatory relationships among genes for understanding regulatory mechanisms in cells. Various methods based on information theory have been developed to infer GRNs. However, these methods introduce many redundant regulatory relationships in the network inference process due to external noise in the original data, topology sparseness in the network structure and non-linear dependency among genes. Especially as the network size increases the performance of these methods decreases dramatically. In this paper, a novel network structure inference method named Loc-PCA-CMI is proposed which first identifies local overlapped gene clusters, and then infers the local network structure for each cluster by a path consistency algorithm based on conditional mutual information (PCA-CMI). The final structure of the GRN is denoted as dependence among genes by ensemble of the obtained local network structures. Loc-PCA-CMI is evaluated on DREAM3 knock-out datasets, and its performance is compared to other information theory based network inference methods including ARACNE, MRNET, PCA-CMI and PCA-PMI. Experimental results demonstrate our novel method Loc-PCA-CMI outperforms other four methods in DREAM 3 datasets especially in size 50 and 100 networks. Source code which includes algorithm and evaluation is freely available for downloading at <http://github.com/chenxofhit/Loc-PCA-CMI>.

**Keywords:** Gene regulatory networks; Network inference; Path consistency algorithm

## 1. Introduction

To infer and understand gene regulatory networks (GRNs) is a critical problem in systems biology, which can help biomedical scientists explicitly to identify complicated regulatory relationships among genes and understand regulatory mechanisms in cells [1,2]. In the past, GRNs were inferred from experimental interventions in which regulatory interactions among genes were verified. Obviously this approach is infeasible [3] and requires substantial time and considerable cost. Owing to the development of micro-array technologies tremendous amounts of gene expression data have been generated [4], which makes it feasible for GRNs to be inferred from these expression data based on computational

methods [5]. In recent years, the inference of networks based on computational methods has become one of the most crucial goals [1,6]. Various methods have been proposed for GRNs inference, such as regression based methods [7–11], ordinary differential equation based methods [12–14], Bayesian and dynamic Bayesian networks [15–20], state-space based methods [21,22]. Unfortunately, gene expression data are typically in high dimensions and relative small sample size which suffer from "dimensionality curse" [23]. Furthermore, gene expression data usually involve large amounts of external noise and non-linear relationships. All of these issues make it more complex and challenging to accurately infer regulatory interactions among genes especially when dealing with large scale gene expression data in the post-genome era.

GRNs can be viewed as undirected acyclic graphs if both up-streaming or down-streaming regulatory relationship among genes are not taken into account and the self-regulatory mechanism is ignored [24], in which each node corresponds to a gene and each edge represents a regulatory relationship between genes. Various computational methods to construct accurate structures of GRNs from expression data have been proposed based on a variety of different assumptions and different conditions [25,26]. Current approaches can be broadly divided into model-based and model-free approaches. Model-based methods usually formulate a computational model of the system and further learn the parameters of such a model. Typical computational models include Boolean network [27–30], Bayesian network [16,31–34], and differential equation models [14,35–39]. Boolean network model is the simplest network model, which is implemented through Boolean variables and Boolean logic. Because the state of gene expression is considered to be only active or inactive, Boolean network models can not entirely capture complex system behavior [40]. The Bayesian network model is a popular probabilistic graphical model in which the dependency relationships among genes are described via a directed acyclic graph (DAG). The Bayesian network model outperforms other models in dealing with noise and incorporating prior knowledge, but structure learning in the model is computationally intensive and has been proved as an NP-hard problem [41]. The differential equation model characterizes the expression level of a gene at a certain time by a function, which involves regulatory interactions with other genes. Differential equation models quantify the change rate (derivative) of the expression of one gene in the system as a function of expression levels of other related genes. A major challenge to use differential equation models for reconstructing GRNs is how to identify the model structure and estimate parameters efficiently in high-dimensional models. Excellent reviews on diverse data-driven modeling schemes and related topics can be found in [42–45].

Other than model-based methods, model-free approaches identify regulatory interactions mainly by measuring dependences between genes. Typical algorithms include correlation-based and information theory-based methods. In the correlation-based method, a regulatory interaction is determined by the degree of co-expression between two genes such as Pearson correlation, rank correlation, Euclidean distance, and the angle between a pair of observed expression vectors [46]. However, the correlation-based methods can not identify complex dependencies between genes, such as non-linear dependencies [47]. Furthermore, quite a few functionally related genes might not be co-expressed, which makes it difficult to accurately infer regulatory interactions. The information theory-based method is also a representative model-free method, in which mutual information (MI) is favored to measure potential dependency among genes as it can capture non-linear dependencies effectively [48,49]. In recent years, various network inference methods based on information theory have been proposed, which focus on distinguishing direct regulatory interactions from indirect associations [50]. To eliminate indirect interactions, Margolin et al. [51] proposed the ARACNE method based on Data Processing Inequality (DPI) with interaction triangles considered. The minimum-redundancy network (MRNET) by Meyer [52] uses a minimum redundancy feature selection method [53], wherein for each candidate gene in a network, it selects a subset of its highly relevant genes while minimizing the MI-based criteria between the selected genes. Zhang et al. [49] introduce a path consistency algorithm based on conditional mutual information (PCA-CMI) and later Zhao et al. [54] introduce a path consistency algorithm based on part mutual information (PCA-PMI). Path consistency algorithm (PCA) is an exhaustive algorithm which is widely

used in inferring GRNs [49]. A trade-off is usually made between running time and accuracy in both PCA-CMI and PCA-PMI. As the network size increases more uncontrollable external noise to the instinct complex network structure makes prediction accuracy of GRNs decrease dramatically. To improve this situation, motivated by the divide and conquer strategy we first use top ranked highly co-expressed genes as centroids of local clusters and then each cluster's accurate structure is refined with PCA-CMI. The final structure of the GRN is then inferred with an ensemble of all the local network structure together. We name this novel approach as Loc-PCA-CMI hereafter and intuitively Loc-PCA-CMI method can deal with relatively larger datasets and benefit from the relatively accurate structure inference for small size gene subnetworks with PCA-CMI.

## 2. Methods

In this section, we introduce related work of information theory including entropy, MI and CMI, as well as the algorithm of Loc-PCA-CMI for inferring GRNs.

### 2.1. Related Work

With the advantages of measuring non-linear dependence association between two variables and relatively high efficiency, information theory is increasingly used to measure the regulatory strength between genes. The mutual information (MI) and conditional mutual information (CMI) are defined as follows:

$$MI(X, Y) = \int \int p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dx dy \quad (1)$$

$$CMI(X, Y|Z) = \int \int \int p(x, y, z) \log \frac{p(x, y|z)}{p(x|z)p(y|z)} \quad (2)$$

where  $p(x, y)$  denotes the joint distribution of two variables (genes)  $X$  and  $Y$ .  $p(x)$  and  $p(y)$  represent the marginal distribution of  $x$  and  $y$ , respectively. For gene expression data, variable  $X$  is a vector, in which the elements denote its expression values in different conditions (samples).

CMI can also be expressed in terms of entropies as :

$$CMI(X, Y|Z) = H(X, Z) + H(Y, Z) - H(Z) - H(X, Y, Z) \quad (3)$$

where  $H(X, Z)$ ,  $H(Y, Z)$ ,  $H(X, Y, Z)$  are joint entropies. High CMI indicates that there may be a close relationship between the variables  $X$  and  $Y$  given variable(s)  $Z$ .

The entropy is estimated with Gaussian kernel probability density estimator [2] and we can get the entropy of variable  $X$  as follows, where  $|C|$  is the determinant of covariance matrix of variable  $X$  [49].

$$H(X) = \log(2\pi e)^{\frac{n}{2}} |C|^{-\frac{1}{2}} \quad (4)$$

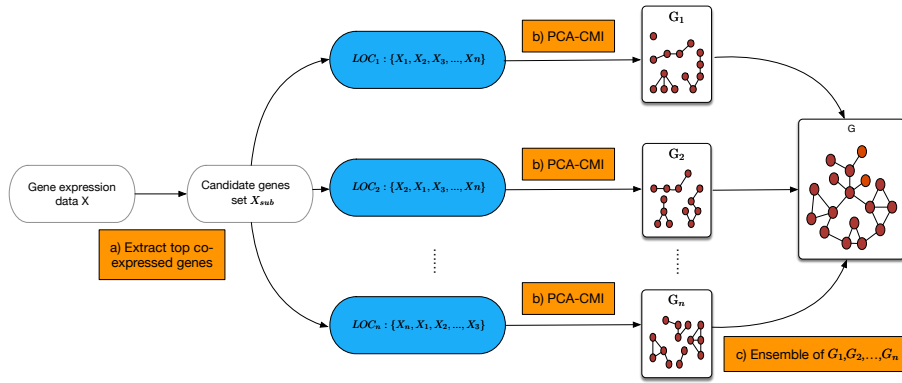
Furthermore, we can obtain the following equation:

$$MI(X, Y) = \frac{1}{2} \log \frac{|C(X)| * |C(Y)|}{|C(X, Y)|} \quad (5)$$

### 2.2. Loc-PCA-CMI

It is well known that biological systems are seldom fully connected and most nodes are only directly connected to a small number of other nodes [55], consequently the GRN is a sparse network. A key step of identifying the sparse structure of the network is to identify the significant edges that may have a comparatively high co-expression value. Specifically our proposed method Loc-PCA-CMI firstly selects the top  $n$  highly co-expressed edges by Pearson correlation analysis with FDR correction

in p-value, and secondly in the reduced edges space computes local overlapped clusters with genes connected by edges. Then for each local cluster we apply the PCA-CMI algorithm, which can construct a high-confidence undirected network [56] by removing the most likely uncorrelated edges repeatedly from low to high order dependence correlation until no edges can be removed, to obtain each local network structure. Final edge weight of the complete regulatory network is obtained by averaging edge weight with each inferred network structure. The entire framework is provided in Figure 1 and the implementation details are shown below in Algorithm 1. As PCA-CMI is extremely competent for relative small GRN structure inference, we make a preprocessing to check the number of gene  $m$  in the algorithm, if  $m$  is less or equal than a constant  $c$  then PCA-CMI is directly applied for the GRN structure inference.



**Figure 1.** The Loc-PCA-CMI framework. a) We first extract top co-expressed genes from gene expression data matrix  $X$  as candidate genes  $X_{sub}$ . The candidate genes  $X_{sub}$  are then grouped into different local clusters with each gene in  $X_{sub}$  as the centroid. b) For each local overlapped cluster PCA-CMI is applied to get accurate structure. c) Ensemble of diverse cluster structure  $G_1, G_2, \dots, G_n$  to obtain the final structure of the GRN as  $G$ .

---

#### Algorithm 1 Loc-PCA-CMI

---

**Input:**  $X$  (gene expression data matrix),  $m$  (number of gene),  $n$  (number of top ranked edges),  $c$  (constant number);  $k$  (CMI order number) and  $\beta$  (order threshold) in subroutine PCA-CMI.

**Output:** Graph weight matrix  $G$

```

1: if  $m \leq c$  then
2:    $G \leftarrow \text{PCA-CMI}(X, k, \beta)$ ;
3:   return  $G$ 
4: else
5:   Construct pair-wise Pearson correlation matrix  $\Omega = \rho(X_i, X_j)$ ;
6:   Select top  $n$  edges as  $E$  with highest Pearson correlation value in  $\Omega$  with FDR correction in
   p-value, and according to which to get the candidate genes as  $X_{sub}$  ;
7:   for each gene  $i$  in  $X_{sub}$  do
8:     Retrieve its directly connected genes that in edges list  $E$  as local cluster  $Loc(i)$ ;
9:     for each cluster  $C$  in  $Loc$  do
10:       $g_C \leftarrow \text{PCA-CMI}(C, k, \beta)$ ;
11:       $G \leftarrow \text{mean}(g_1, g_2, \dots, g_t)$ , where  $t$  stands for total cluster number of  $Loc$ ;
12:   return  $G$ 

```

---

The computational complexity of Algorithm 1 is generally determined by total cluster number  $t$  and the complexity of the subroutine PCA-CMI.  $t$  is in the same order with the number of genes  $m$ . The computational complexity of PCA-CMI are dominated by the parameter CMI order number  $k$  and cluster size  $|C|$  of  $C$ , which can be roughly estimated as  $O(|C|^k)$ . As a result the final computational

complexity of PCA-CMI can be calculated as  $O(m * |C|^k)$ . At worst, if cluster size  $|C|$  equals to  $m$  i.e. every cluster contains all the genes in it and the computational complexity is thus  $m * m^k = m^{k+1}$ . However, this worst case scenario seldom happens in practice; and actually  $|C|$  is usually much lower than  $m$ .

### 3. Materials

We benchmarked the performance of our approach, Loc-PCA-CMI using six simulation data from well known DREAM3 challenge [57]. DREAM3 features in silico networks and expression data simulated using GeneNetWeaver software. Benchmark networks were derived as subnetworks of a system of regulatory interactions from known model organisms: E.coli and S.cerevisiae. Six gene knock-out expression networks in DREAM3 are evaluated in our experiments, which include three different sizes varying in 10, 50, 100 with two types of model organisms E.coli and S.cerevisiae respectively. Table 1 gives detailed descriptions of the datasets.

**Table 1.** Descriptions of the datasets in our experiments

Datasets	#Samples	#Average(Max) degree	#Edges	#Network density
DREAM3-10 Ecoli	11	2.2(5)	11	0.244
DREAM3-50 Ecoli	51	2.48(14)	62	0.051
DREAM3-100 Ecoli	101	2.5(14)	125	0.025
DREAM3-10 Yeast	11	2(4)	10	0.222
DREAM3-50 Yeast	51	3.08(13)	77	0.063
DREAM3-100 Yeast	101	3.32(10)	166	0.034

Taken the input datafile of DREAM3-10 Ecoli as an example, lines stand for samples and 10 columns stand for 10 different gene expression data. The first line is the wild-type expression data, every gene in this sample stays at a steady state. The  $k$  th ( $2 \leq k \leq 11$ ) line stands for that how other gene expression data changes after the  $k - 1$  th gene is knock-out.

### 4. Results and Discussion

As described in Algorithm 1 three intrinsic parameters affect the performance of Loc-PCA-CMI in GRN structure inference. The first parameter is the number  $n$  of top selected edges. If  $n$  increases more edges are considered and the local cluster size will increase subsequently. The second parameter is  $\beta$  which acts as the threshold value of MI and CMI to decide independence. The third parameter is CMI order number  $k$ , theoretically by increasing  $k$  the structure is more accurate if CMI does not reach the threshold  $\beta$  in  $k - 1$  order. Latter two parameters are with PCA-CMI and PCA-PMI. The best value of  $n$  can be obtained by cross validation and generally the larger value of  $n$  can contribute to a larger size cluster and more genes are covered in the network, in our experiments we set its value to be  $n = 20\% * \binom{m}{2}$  uniformly. Besides the above three intrinsic parameters we set constant  $c = 10$  in Algorithm 1, i.e. if the number of genes is less or equal to 10 Loc-PCA-CMI calls PCA-CMI directly and in this case performance of Loc-PCA-CMI and PCA-CMI are the same.

We assess the performance of Loc-PCA-CMI by evaluating the area under the Receiver Operating Characteristic curve (AUROC) and the area under the Precision-Recall curve (AUPR). As in sparse biological networks the number of non-existing edges (negatives) outweighs the number of existing edges (positives) significantly, AUPR is more informative to AUROC in fact [58]. We tend to use AUPR for evaluation, but for a conservative comparison with other methods that adopt AUROC as evaluation metric we also take AUROC as a supplementary metric. Higher AUROC and AUPR value indicate more accurate GRN predictions. For this purpose, we compute the numbers of true positive (TP), true negative (TN), false positive (FP) and false negative (FN) edges by comparing the regulatory edges in the gold standard network with the top  $q$  edges from the ranked list output of Loc-PCA-CMI.

The ROC curve is constructed by plotting the true positive rates ( $\text{TPR} = \text{TP}/(\text{TP}+\text{FN})$ ) versus the false positive rates ( $\text{FPR} = \text{FP}/(\text{FP}+\text{TN})$ ) for increasing  $q$  ( $q = 1, 2, \dots, m^2$ ). Similarly, the precision ( $\text{TP}/(\text{TP}+\text{FP})$ ) and recall ( $\text{TP}/(\text{TP}+\text{FN})$ ) curve is plotted for increasing  $q$ .

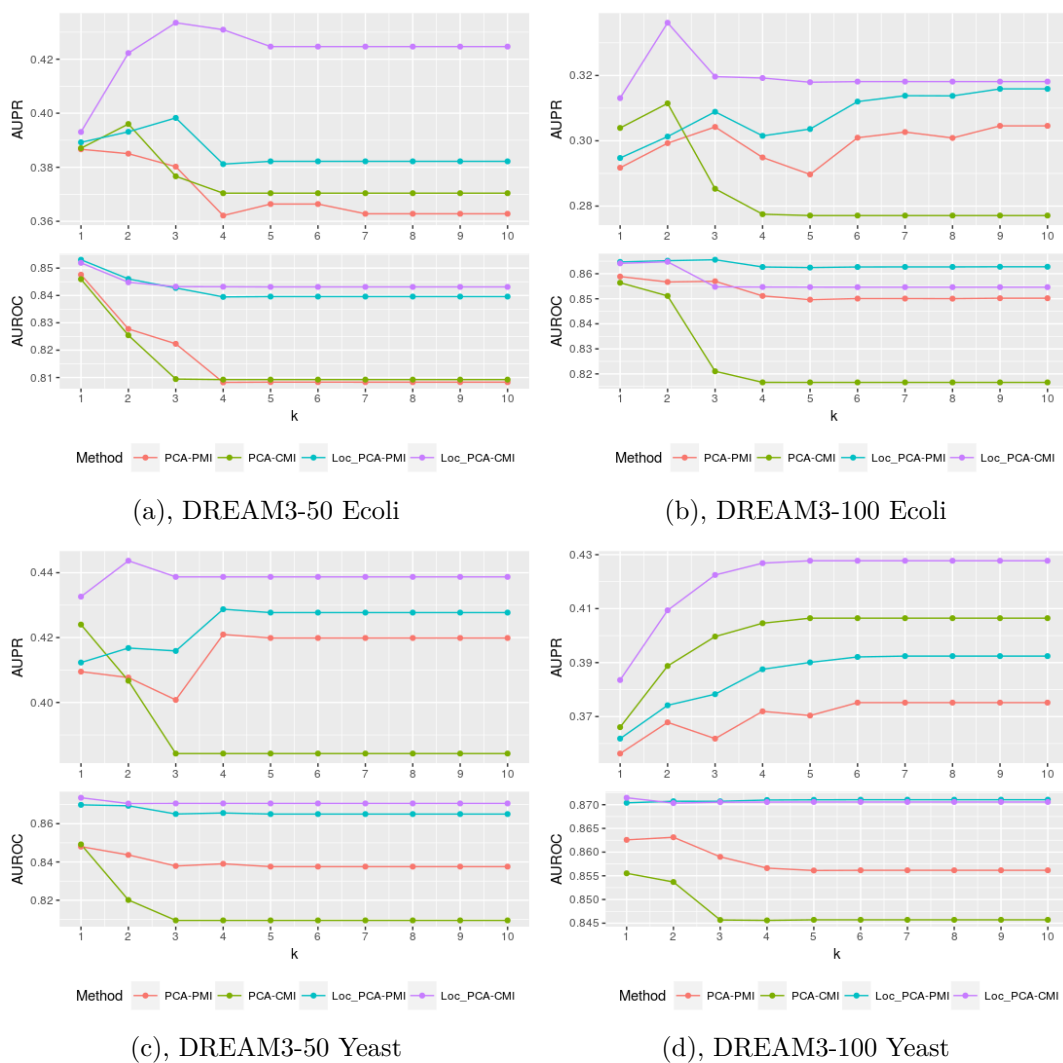
It should be noticed that in Algorithm 1 after each local cluster is obtained both PCA-CMI and PCA-PMI are alternatives for the subsequent structure refinement. If PCA-CMI is replaced with PCA-PMI a novel method is generated and we name it as Loc-PCA-PMI analogously. Then four PCA based methods are derived including PCA-PMI, PCA-CMI, Loc-PCA-PMI, Loc-PCA-CMI at present, all of which belong to model-free methods. As showed in Table 1 among the six benchmark datasets DREAM10-Ecoli and DREAM10-Yeast datasets contain only 10 genes, hence Loc-PCA-CMI and PCA-CMI are identical in performance and the same with Loc-PCA-PMI and PCA-PMI according to the principle of Algorithm 1. For meaningful comparison of these PCA based methods we select other four datasets whose gene number is greater than 10. Order number is not explicitly discussed in [49,54] wherein  $\beta = 0.03$  and  $k = 2$  are set directly, we are curious about how order number  $k$  affects the performance of these methods as well. By varying the order number  $k$  from 1 to 10 in these four methods with fixed threshold  $\beta = 0.03$ , AUROC and AUPR can be calculated, respectively. Figure 2 illustrates the result in summary on the benchmark datasets, and from which we can draw three conclusions:

- Order number  $k$  affects the results of these four PCA based methods, generally when  $k$  reaches 4 AUPR and AUROC become stable except those in DREAM3-100 Ecoli dataset.
- Loc-PCA-CMI and Loc-PCA-PMI yield higher AUPR and AUROC than PCA-CMI and PCA-PMI, respectively, hence the local cluster strategy adopted in the algorithm helps to improve the performance of PCA-CMI and PCA-PMI.
- Loc-PCA-CMI outperforms than the other three methods in the metric of AUPR, which is more meaningful for sparse network structure prediction issues.

We also conduct a comparison experiment using Loc-PCA-CMI with four previously proposed methods on all the six benchmark datasets, which include ARACNE, MRNET, PCA-PMI, PCA-CMI. We use the R package "minet" with default parameters for evaluation of ARACNE and MRNET [59]. The MI matrices of the methods are approximated by using Pearson correlation directly from continuous gene knock-out expression data [60,61]. For implementation of PCA-PMI and PCA-CMI we have downloaded the MATLAB codes according to the URL provided in [49,54]. We prefer the default value of parameters in PCA-PMI and PCA-CMI, where  $\beta = 0.03$  and  $k = 2$ . For Loc-PCA-CMI we also adopt the same values to these two parameters for comparison. Table 2 gives the AUROC and AUPR of this experiment. From the table we can see that AUPR of all these contending methods decrease dramatically when the network size increases. Loc-PCA-CMI is only after PCA-PMI in DREAM3-10 Yeast dataset, while in other five datasets it outperforms the other four methods ARACNE, MRNET, PCA-PMI, PCA-CMI in terms of both AUROC and AUPR. We provide the source codes including all the methods, benchmark datasets and evaluation scripts at <http://github.com/chenxofhit/Loc-PCA-CMI>.

**Table 2.** AUROC and AUPR for the six datasets using different methods

Dataset	ARACNE		MRNET		PCA-PMI		PCA-CMI		Loc-PCA-CMI	
	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
DREAM3-10 Ecoli	0.523	0.255	0.518	0.258	0.816	0.483	0.825	0.499	<b>0.825</b>	<b>0.499</b>
DREAM3-50 Ecoli	0.474	0.050	0.529	0.061	0.828	0.385	0.825	0.396	<b>0.845</b>	<b>0.422</b>
DREAM3-100 Ecoli	0.505	0.027	0.488	0.025	0.857	0.299	0.851	0.311	<b>0.865</b>	<b>0.336</b>
DREAM3-10 Yeast	0.628	0.321	0.644	0.322	<b>0.995</b>	<b>0.933</b>	0.993	0.918	0.993	0.918
DREAM3-50 Yeast	0.507	0.074	0.524	0.080	0.844	0.408	0.820	0.406	<b>0.871</b>	<b>0.444</b>
DREAM3-100 Yeast	0.547	0.040	0.556	0.042	0.863	0.368	0.854	0.389	<b>0.870</b>	<b>0.409</b>



**Figure 2.** AUPR and AUROC by varying  $k$  from 1 to 10 of four PCA based methods on four different datasets: (a) DREAM3-50 Ecoli; (b) DREAM3-100 Ecoli; (c) DREAM3-50 Yeast; (d) DREAM3-100 Yeast.



## 5. Conclusion

We have proposed a novel model-free gene regulatory network structure inference method named Loc-PCA-CMI, which is motivated by the divide and conquer strategy. At present all the experiments are conducted in the DREAM 3 challenge silico datasets. Experiments on DREAM3 knock-out datasets show that Loc-PCA-CMI benefits from the local overlapped cluster strategy. Besides, Loc-PCA-CMI outperforms other comparing methods including ARACNE, MRNET, PCA-PMI, PCA-CMI especially for networks with sizes of 50 and 100.

Loc-PCA-CMI is an extended version of PCA-CMI and its limitation in computational inefficiency can be inherited, especially when to deal with large size datasets. The number of local clusters in case of large network can be extremely large. However if we can control the size of each local cluster, our method should be applicable on large size data. One of our future works is to improve the cluster strategy to be more efficient and effective to deal with large size data. We mainly focus on inferring GRNs structure and have not considered the stability of networks in this study. As a result another one of our future works is to infer stable networks.

**Acknowledgments:** This work was supported in part by the National Natural Science Foundation of China under Grants (No. 61622213, No. 61732009 and No. 61772552). We would also like to thank four anonymous reviewers for comments and suggestions to improve this manuscript.

**Author Contributions:** Xiang Chen contributed to the design of the study and experiments, invention of the Loc-PCA-CMI method, data analysis, and writing of the paper. Min Li, Ruiqing Zheng and Siyu Zhao proposed constructive suggestion about the detail of the experiments. Min Li and Fang-Xiang Wu reviewed the paper. All authors read and approved the final manuscript.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

- Altay, G.; Emmert-Streib, F. Inferring the conservative causal core of gene regulatory networks. *BMC Systems Biology* **2010**, *4*, 132.
- Basso, K.; Margolin, A.A.; Stolovitzky, G.; Klein, U.; Dalla-Favera, R.; Califano, A. Reverse engineering of regulatory networks in human B cells. *Nature genetics* **2005**, *37*, 382.
- Elnitski, L.; Jin, V.X.; Farnham, P.J.; Jones, S.J. Locating mammalian transcription factor binding sites: a survey of computational and experimental techniques. *Genome research* **2006**, *16*, 1455–1464.
- Hughes, T.R.; Marton, M.J.; Jones, A.R.; Roberts, C.J.; Stoughton, R.; Armour, C.D.; Bennett, H.A.; Coffey, E.; Dai, H.; He, Y.D.; others. Functional discovery via a compendium of expression profiles. *Cell* **2000**, *102*, 109–126.
- Maetschke, S.R.; Madhamshettiwar, P.B.; Davis, M.J.; Ragan, M.A. Supervised, semi-supervised and unsupervised inference of gene regulatory networks. *Briefings in bioinformatics* **2013**, *15*, 195–211.
- Margolin, A.A.; Wang, K.; Lim, W.K.; Kustagi, M.; Nemenman, I.; Califano, A. Reverse engineering cellular networks. *Nature protocols* **2006**, *1*, 662.
- Huynh-Thu, V.A.; Irrthum, A.; Wehenkel, L.; Geurts, P. Inferring regulatory networks from expression data using tree-based methods. *PLoS One* **2010**, *5*, 1–10.
- Haury, A.C.; Mordelet, F.; Vera-Licona, P.; Vert, J.P. TIGRESS: Trustful Inference of Gene REgulation using Stability Selection. *BMC Syst. Biol.* **2012**, *6*, 145.
- Huynh-Thu, V.A.; Sanguinetti, G.; Huynh-thu, A.; Jump, T. Combining tree-based and dynamical systems for the inference of gene regulatory networks. *Bioinformatics* **2014**, *31*, 1614–1622.
- Liu, L.Z.; Wu, F.X.; Zhang, W.J. A group LASSO-based method for robustly inferring gene regulatory networks from multiple time-course datasets. *BMC systems biology* **2014**, *8*, S1.
- Li, M.; Zheng, R.; Li, Y.; Wu, F.X.; Wang, J. MGT-SM: A Method for Constructing Cellular Signal Transduction Networks. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **2017**.
- Sakamoto, E.; Iba, H. Inferring a system of differential equations for a gene regulatory network by using genetic programming. *Evolutionary Computation*, 2001. Proceedings of the 2001 Congress on. IEEE, 2001, Vol. 1, pp. 720–726.



13. Chowdhury, A.R.; Chetty, M.; Evans, R. Stochastic S-system modeling of gene regulatory network. *Cognitive neurodynamics* **2015**, *9*, 535–547.
14. Li, Z.; Li, P.; Krishnan, A.; Liu, J. Large-scale dynamic gene regulatory network inference combining differential equation models with local dynamic Bayesian network analysis. *Bioinformatics* **2011**, *27*, 2686–2691.
15. Murphy, K.; Mian, S.; others. Modelling gene expression data using dynamic Bayesian networks. Technical report, Technical report, Computer Science Division, University of California, Berkeley, CA, 1999.
16. Zou, M.; Conzen, S.D. A new dynamic Bayesian network (DBN) approach for identifying gene regulatory networks from time course microarray data. *Bioinformatics* **2004**, *21*, 71–79.
17. Vinh, N.X.; Chetty, M.; Coppel, R.; Wangikar, P.P. GlobalMIT: learning globally optimal dynamic bayesian network with the mutual information test criterion. *Bioinformatics* **2011**, *27*, 2765–2766.
18. Young, W.C.; Raftery, A.E.; Yeung, K.Y. Fast Bayesian inference for gene regulatory networks using ScanBMA. *BMC systems biology* **2014**, *8*, 47.
19. Liu, F.; Zhang, S.W.; Guo, W.F.; Wei, Z.G.; Chen, L. Inference of Gene Regulatory Network Based on Local Bayesian Networks. *PLOS Comput. Biol.* **2016**, *12*, e1005024.
20. Omranian, N.; Eloundou-Mbebi, J.M.; Mueller-Roeber, B.; Nikoloski, Z. Gene regulatory network inference using fused LASSO on multiple data sets. *Scientific reports* **2016**, *6*, 20533.
21. Wu, F.X.; Zhang, W.J.; Kusalik, A.J. Modeling gene expression from microarray expression data with state-space equations. In *Biocomputing 2004*; World Scientific, 2003; pp. 581–592.
22. Quach, M.; Brunel, N.; d'Alché Buc, F. Estimating parameters and hidden variables in non-linear state-space models based on ODEs for biological networks inference. *Bioinformatics* **2007**, *23*, 3209–3216.
23. Wang, Y.; Joshi, T.; Zhang, X.S.; Xu, D.; Chen, L. Inferring gene regulatory networks from multiple microarray datasets. *Bioinformatics* **2006**, *22*, 2413–2420.
24. Irrthum, A.; Wehenkel, L.; Geurts, P.; others. Inferring regulatory networks from expression data using tree-based methods. *PloS one* **2010**, *5*, e12776.
25. Longabaugh, W.J.; Davidson, E.H.; Bolouri, H. Computational representation of developmental genetic regulatory networks. *Developmental biology* **2005**, *283*, 1–16.
26. Karlebach, G.; Shamir, R. Modelling and analysis of gene regulatory networks. *Nature reviews. Molecular cell biology* **2008**, *9*, 770.
27. Shmulevich, I.; Dougherty, E.R.; Kim, S.; Zhang, W. Probabilistic Boolean networks: a rule-based uncertainty model for gene regulatory networks. *Bioinformatics* **2002**, *18*, 261–274.
28. Kim, H.; Lee, J.K.; Park, T. Boolean networks using the chi-square test for inferring large-scale gene regulatory networks. *BMC bioinformatics* **2007**, *8*, 37.
29. Bornholdt, S. Boolean network models of cellular regulation: prospects and limitations. *Journal of the Royal Society Interface* **2008**, *5*, S85–S94.
30. Zhou, J.X.; Samal, A.; d'Hérouël, A.F.; Price, N.D.; Huang, S. Relative stability of network states in Boolean network models of gene regulation in development. *Biosystems* **2016**, *142*, 15–24.
31. Kim, S.Y.; Imoto, S.; Miyano, S. Inferring gene networks from time series microarray data using dynamic Bayesian networks. *Briefings in bioinformatics* **2003**, *4*, 228–235.
32. Chen, X.w.; Anantha, G.; Wang, X. An effective structure learning method for constructing gene networks. *Bioinformatics* **2006**, *22*, 1367–1374.
33. Needham, C.J.; Bradford, J.R.; Bulpitt, A.J.; Westhead, D.R. A primer on learning in Bayesian networks for computational biology. *PLoS computational biology* **2007**, *3*, e129.
34. Lo, L.Y.; Wong, M.L.; Lee, K.H.; Leung, K.S. High-order dynamic Bayesian Network learning with hidden common causes for causal gene regulatory network. *BMC bioinformatics* **2015**, *16*, 395.
35. Gardner, T.S.; Di Bernardo, D.; Lorenz, D.; Collins, J.J. Inferring genetic networks and identifying compound mode of action via expression profiling. *Science* **2003**, *301*, 102–105.
36. di Bernardo, D.; Thompson, M.J.; Gardner, T.S.; Chobot, S.E.; Eastwood, E.L.; Wojtovich, A.P.; Elliott, S.J.; Schaus, S.E.; Collins, J.J. Chemogenomic profiling on a genome-wide scale using reverse-engineered gene networks. *Nature biotechnology* **2005**, *23*, 377–383.
37. Bansal, M.; Gatta, G.D.; Di Bernardo, D. Inference of gene regulatory networks and compound mode of action from time course gene expression profiles. *Bioinformatics* **2006**, *22*, 815–822.

38. Honkela, A.; Girardot, C.; Gustafson, E.H.; Liu, Y.H.; Furlong, E.E.; Lawrence, N.D.; Rattray, M. Model-based method for transcription factor target identification with limited data. *Proceedings of the National Academy of Sciences* **2010**, *107*, 7793–7798.
39. Lu, T.; Liang, H.; Li, H.; Wu, H. High-dimensional ODEs coupled with mixed-effects modeling techniques for dynamic gene regulatory network identification. *Journal of the American Statistical Association* **2011**, *106*, 1242–1258.
40. Lee, W.P.; Tzou, W.S. Computational methods for discovering gene networks from expression data. *Briefings in bioinformatics* **2009**, *10*, 408–423.
41. Chickering, D.M.; Heckerman, D.; Meek, C. Large-sample learning of Bayesian networks is NP-hard. *Journal of Machine Learning Research* **2004**, *5*, 1287–1330.
42. Hecker, M.; Lambeck, S.; Toepfer, S.; Van Someren, E.; Guthke, R. Gene regulatory network inference: data integration in dynamic models—a review. *Biosystems* **2009**, *96*, 86–103.
43. Marbach, D.; Costello, J.C.; Küffner, R.; Vega, N.M.; Prill, R.J.; Camacho, D.M.; Allison, K.R.; Kellis, M.; Collins, J.J.; Stolovitzky, G.; others. Wisdom of crowds for robust gene network inference. *Nature methods* **2012**, *9*, 796–804.
44. Wu, F.X. Inference of gene regulatory networks and its validation. *Current Bioinformatics* **2007**, *2*, 139–144.
45. Liu, L.Z.; Wu, F.X.; Zhang, W.J. Reverse engineering of gene regulatory networks from biological data. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* **2012**, *2*, 365–385.
46. Wang, Y.R.; Huang, H. Review on statistical methods for gene network reconstruction using expression data. *Journal of theoretical biology* **2014**, *362*, 53–61.
47. Ruysinck, J.; Geurts, P.; Dhaene, T.; Demeester, P.; Saeys, Y.; others. Nimefi: gene regulatory network inference using multiple ensemble feature importance algorithms. *PLoS One* **2014**, *9*, e92709.
48. Brunel, H.; Gallardo-Chacón, J.J.; Buil, A.; Vallverdú, M.; Soria, J.M.; Caminal, P.; Perera, A. MISS: a non-linear methodology based on mutual information for genetic association studies in both population and sib-pairs analysis. *Bioinformatics* **2010**, *26*, 1811–1818.
49. Zhang, X.; Zhao, X.M.; He, K.; Lu, L.; Cao, Y.; Liu, J.; Hao, J.K.; Liu, Z.P.; Chen, L. Inferring gene regulatory networks from gene expression data by path consistency algorithm based on conditional mutual information. *Bioinformatics* **2011**, *28*, 98–104.
50. Marbach, D.; Prill, R.J.; Schaffter, T.; Mattiussi, C.; Floreano, D.; Stolovitzky, G. Revealing strengths and weaknesses of methods for gene network inference. *Proceedings of the national academy of sciences* **2010**, *107*, 6286–6291.
51. Margolin, A.A.; Nemenman, I.; Basso, K.; Wiggins, C.; Stolovitzky, G.; Dalla Favera, R.; Califano, A. ARACNE: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC bioinformatics* **2006**, *7*, S7.
52. Meyer, P.E.; Kontos, K.; Lafitte, F.; Bontempi, G. Information-theoretic inference of large transcriptional regulatory networks. *EURASIP journal on bioinformatics and systems biology* **2007**, *2007*, 79879.
53. Peng, H.; Long, F.; Ding, C. Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on pattern analysis and machine intelligence* **2005**, *27*, 1226–1238.
54. Zhao, J.; Zhou, Y.; Zhang, X.; Chen, L. Part mutual information for quantifying direct associations in networks. *Proceedings of the National Academy of Sciences* **2016**, *113*, 5130–5135.
55. Jeong, H.; Tombor, B.; Albert, R.; Oltvai, Z.N.; Barabási, A.L. The large-scale organization of metabolic networks. *Nature* **2000**, *407*, 651–654.
56. Spirtes, P.; Glymour, C.N.; Scheines, R. *Causation, prediction, and search*; MIT press, 2000.
57. Schaffter, T.; Marbach, D.; Floreano, D. GeneNetWeaver: in silico benchmark generation and performance profiling of network inference methods. *Bioinformatics* **2011**, *27*, 2263–2270.
58. Saito, T.; Rehmsmeier, M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PloS one* **2015**, *10*, e0118432.
59. Meyer, P.E.; Lafitte, F.; Bontempi, G. minet: AR/Bioconductor package for inferring large transcriptional networks using mutual information. *BMC bioinformatics* **2008**, *9*, 461.

- 342 60. Olsen, C.; Meyer, P.E.; Bontempi, G. On the impact of entropy estimation on transcriptional regulatory  
343 network inference based on mutual information. *EURASIP Journal on Bioinformatics and Systems*  
344 *Biology* **2008**, *2009*, 308959.
- 345 61. Meyer, P.; Marbach, D.; Roy, S.; Kellis, M. Information-Theoretic Inference of Gene Networks Using  
346 Backward Elimination. *BioComp*, 2010, pp. 700–705.

347 © 2017 by the authors. Submitted to *Int. J. Mol. Sci.* for possible open access  
348 publication under the terms and conditions of the Creative Commons Attribution (CC BY) license  
349 (<http://creativecommons.org/licenses/by/4.0/>).