

# Gene Expression Level Association Analysis for Yeast Data

## Group 5

### 1. Yeast Data

We explored the effects of different statistical methods by the analysis of a yeast eQTL data set described by Brem and Kruglyak (2005), where  $n = 112$  segregants were grown from a cross between two budding yeast strains, BY4716 and RM11-1a. For each of the segregants, gene expression was profiled on microarrays containing 6216 genes, and genotyping was performed at 2957 markers.

#### 1.1 Data Preparation

##### 1.1.1 PROCESSING MARKERS DATA

The DNA of the same organism has a very high similarity. For example, all humans are 99.9% identical and, of that tiny 0.1% difference. Considering the similarity of markers, we found that different markers in the data may have identical or differ at most few yeast samples after exploratory data analysis.

Inspired by Yin and Li (2011), we combined the markers into 949 blocks such that markers with the same block differed by at most one sample, and one representative marker was chosen from each block. To accomplish that, we used hierarchical clustering with complete-linkage criteria and cut the clustering tree at distance  $d = 1$ . Note that the value of the marker data is only 0 or 1, so we can use any distance function, like Euclidean distance or Manhattan distance, to obtain the same result.

We can further reduce the number of markers, observing that different markers have different functions, and some markers have no effect on the gene expression we want to study. Thus a marginal gene-marker association analysis was then performed to identify markers that are associated with the expression levels of at least two genes with a  $p$ -value less than 0.05, resulting in a total of  $p = 776$  markers with additional 18.2% reduction .

##### 1.1.2 PROCESSING EXPRESSION LEVEL DATA

We choose genes according to MAPK signaling pathways (Kanehisa et al., 2013) from *The yeast MAPK pathway from the KEGG database*<sup>1</sup> as shown in Figure 1.

From the database we obtain 53 markers which have been linked to MAPK signaling pathways biologically.

After the above operations, we got the processed data with a  $X \in \mathbb{R}^{112 \times 949}$ ,  $Y \in \mathbb{R}^{112 \times 53}$ .

---

1. <http://www.genome.jp/kegg/pathway/sce/sce04011.html>

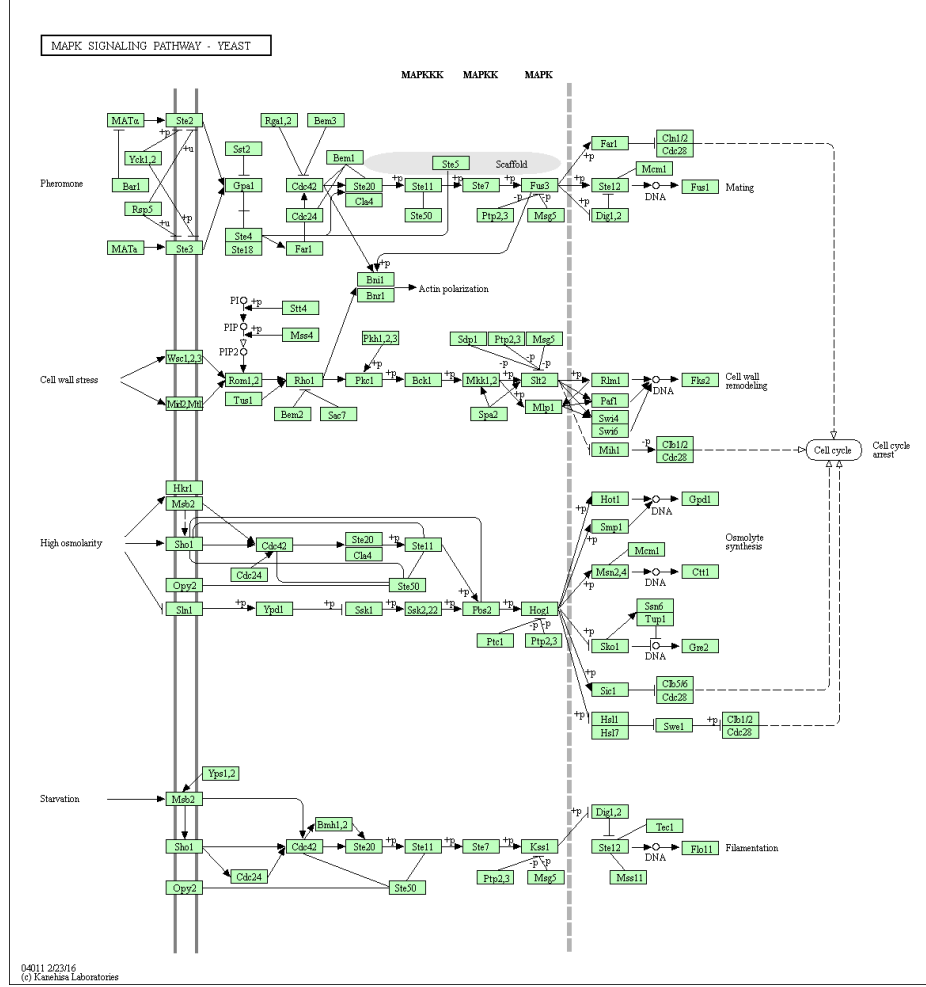


Figure 1: The yeast MAPK pathway from the KEGG database

Hereinafter, we denote  $p$  as the number of explanatory variables,  $q$  as the number of response variable,  $n$  as the number of samples,  $E$  as the random error matrix, and  $B$  as the coefficient matrix, we can construct a multi-response linear model as following

$$Y = XB + E. \quad (1)$$

## 2. Methodology

Here we are going to introduce methods we used to select variables.

As described in (1), we are facing a linear regression model in high dimensional space with multiple response variables. This makes it more difficult for us to deal with problems.

Next, we will introduce three methods to solve the model. They were verified as a good method to deal with univariate high-dimensional linear regression, grouped variables, and factorization problem.

## 2.1 Unit-Response Regression

Denote  $Y_j$  as the  $j$ -th column of  $Y$ , which represents the expression level of  $j$ -th gene.

It is intuitive for us to divide and conquer it by regressing each  $Y_j$  with  $X$  and we will get  $q$  linear models.

We can use classical LASSO (Tibshirani, 1996) to get the estimated coefficient vector  $\hat{\beta}_{\cdot j} \in \mathbb{R}^p$ . Then combine  $q$  coefficient vectors  $\hat{\beta}_{\cdot j}$  into a matrix  $\hat{B} \in \mathbb{R}^{p \times q}$  by column.

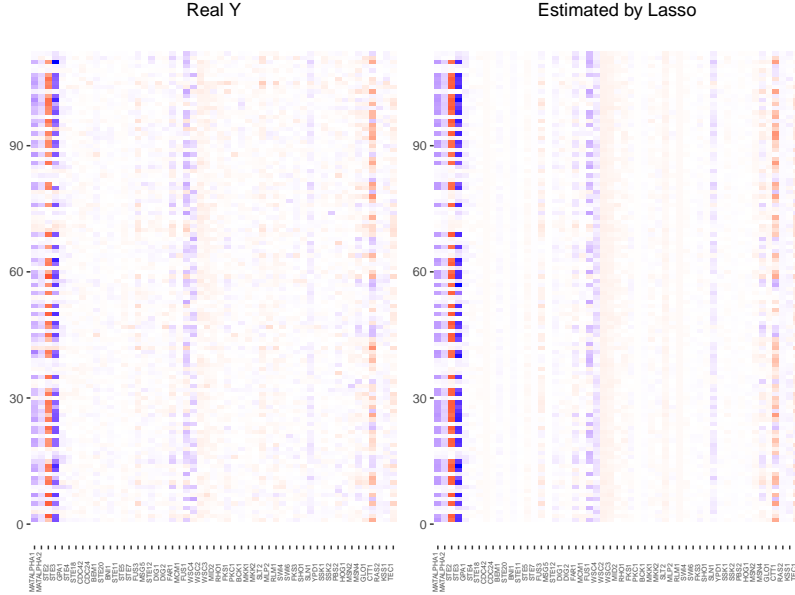


Figure 2: Heatmaps of real  $Y$  and  $\hat{Y}$  by LASSO

We plot heatmap of predicted  $Y$  and the real one in Figure 2 because it is simpler and more intuitive than displaying numbers. We can learn from the diagram that the two colors are almost the same, indicating that the estimation results by LASSO are acceptable.

However, as is known, LASSO is a kind of shrinkage estimation method with  $L_1$  penalty, which not only shrinks the size of the coefficients, but also set some of them to zero. So a sparse  $\hat{\beta}_{\cdot j}$  is expected for each  $j \in \{1, 2, \dots, q\}$ . We rewrite the optimization function as following

$$\hat{B} = \arg \min_B \left\{ \frac{1}{2n} \|Y - XB\|_F^2 + \sum_{i=1}^q \lambda_i \|B_{\cdot i}\|_1 \right\} \quad (2)$$

where  $\|\cdot\|_F$  is the Frobenius norm and  $\sum_{i=1}^q \lambda_i \|B_{\cdot i}\|_1$  can be regraded as a  $(1, 1)$  norm of matrix  $B$ .

Actually, we got 602 nonzero rows in  $\hat{B}$  which also has full column rank.

## 2.2 Grouped-Response Regression

Group sparse linear regression for multitask learning (Dai and Barber, 2016)

$$\hat{B} = \arg \min_B \left\{ \frac{1}{2n} \|Y - XB\|_F^2 + \lambda \|B\|_{(2,1)} \right\} \quad (3)$$

where the  $(2, 1)$  norm in the penalty is given by  $\|B\|_{(2,1)} = \sum_i \sqrt{\sum_j B_{ij}^2}$ .

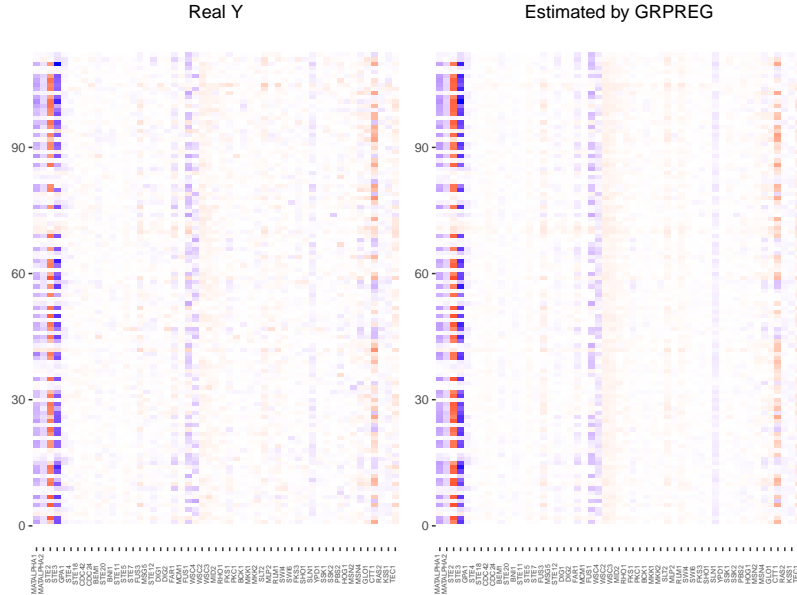


Figure 3: Heatmaps.  $\hat{B}$  has full column rank and 201 non-zero rows.

As shown in the Figure 3, although it totally uses fewer predictors variables than LASSO, but it also restores  $Y$  well.

Intuitively,  $(2, 1)$  norm penalty can be understood as a  $L_1$  penalty imposed on sum of squares of the rows of  $B$ . If we choose a large  $\lambda$ ,  $\hat{B}$  will contain many zero rows, however the nonzero rows will themselves be dense (not entrywise sparsity). This is because  $(2, 1)$  norm penalty promotes rowwise sparsity of  $\hat{B}$ .

Biologically, the result shows that among all  $p$  markers, a total of 201 markers are related to 53 genes. And each marker selected has some effect on each gene in our data.

### 2.3 Multi-Response Regression

Previous biological research has revealed some facts which could be the guidelines for us to choose suitable method of data analysis.

Extensive genetic and biochemical analysis has revealed that there are a few functionally distinct signaling pathways of genes (Gustin et al., 1998; Brem and Kruglyak, 2005), suggesting that the association structure between the eQTLs and the genes is of low rank. Each signaling pathway involves only a subset of genes, which are regulated by only a few genetic variants, suggesting that each association between the eQTLs and the genes is sparse in both the input and the output (or in both the responses and the predictors), and the pattern of sparsity should be pathway specific. Moreover, it is known that the yeast MAPK pathways regulate and interact with each other (Gustin et al., 1998).

The complex genetic structures described above clearly indicate that the association structure between the eQTLs and the gene is of low rank and sparsity. It call for a joint

statistical analysis that can reveal multiple distinct associations between subsets of genes and subsets of genetic variants.

SOFAR (Uematsu et al., 2019) uses the SVD decomposition  $B = UDV^T$  and then impose penalties into  $U$ ,  $D$  and  $V$  respectively.

$$\begin{aligned} (\hat{D}, \hat{U}, \hat{V}) = \arg \min_{D, U, V} & \left\{ \frac{1}{2n} \|X - UDV^T\|_F^2 + \lambda_d \|D\|_1 + \lambda_a \rho_a(UD) + \lambda_b \rho_b(VD) \right\} \\ & \text{subject to } U^T U = \mathbf{I}_m, \quad V^T V = \mathbf{I}_m \end{aligned} \quad (4)$$

Rank reduction is achieved mainly through the first term and variable selection is achieved through the last two terms.  $\rho_a(\cdot)$  and  $\rho_b(\cdot)$  can be equal or distinct, depending on our questions and goals.  $\rho$  can be entrywise  $L_1$  norm or rowwise  $(2, 1)$  norm or others.

Figure 4 shows the results from SOFAR with turning parameters  $\lambda_d$ ,  $\lambda_a$  and  $\lambda_b$  choosen by cross validation. SOFAR returned rowwise sparse singular matrix  $U$  and  $V$  and a low rank diagram matrix  $D$ . It finds only 228 significant markers by selecting non-zero rows in the estimation of  $U$ , 25 genes by selecting non-zero rows in the estimation of  $V$  and 3 latent patterns within the data.

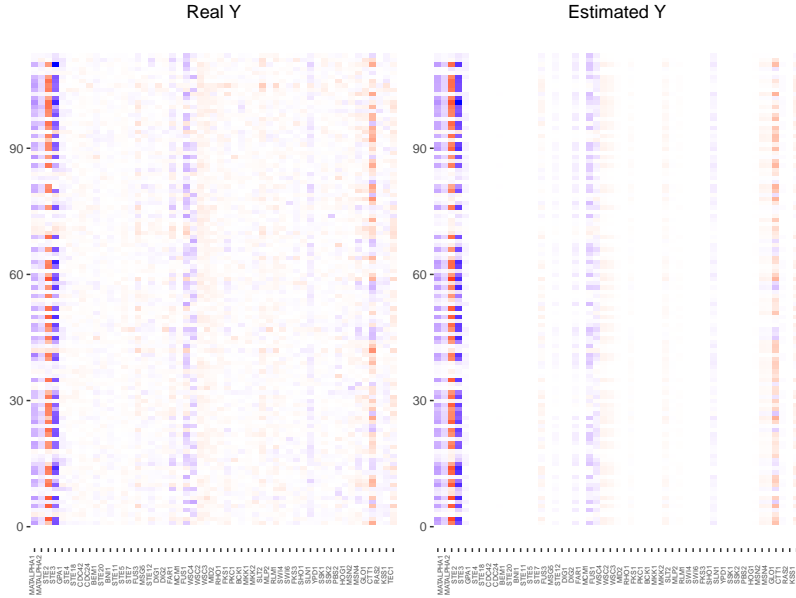


Figure 4: Heatmaps of real  $Y$  and  $\hat{Y}$  by SOFAR

Specifically, the rank of coefficient matrix is 3, as shown in the Figure 5, indicating that there are only three latent variables ( $XU$ ). Figure 5 shows the specific linear relationship between  $XU$  and response variables  $YV$  in different patterns, and the fitting is pretty good. And here we also get the sparsity of  $U$  and  $V$ .

In order to check whether the SOFAR method completely finds the latent pattern within the data, we drew the fourth possible latent pattern in Figure 6. As is shown, the latent response and the latent predictor are almost completely independent and there is no possible relationship.

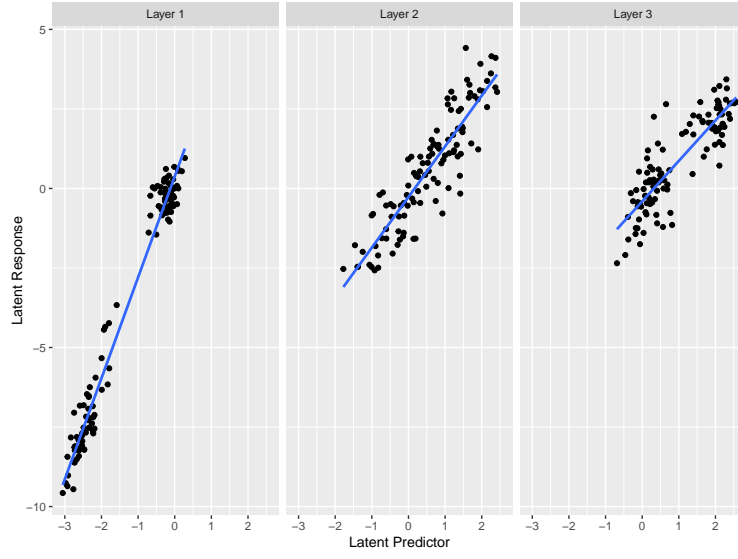


Figure 5: Scatter plots of the latent responses versus the latent predictors in three SVD layers for the yeast data estimated by the SOFAR method

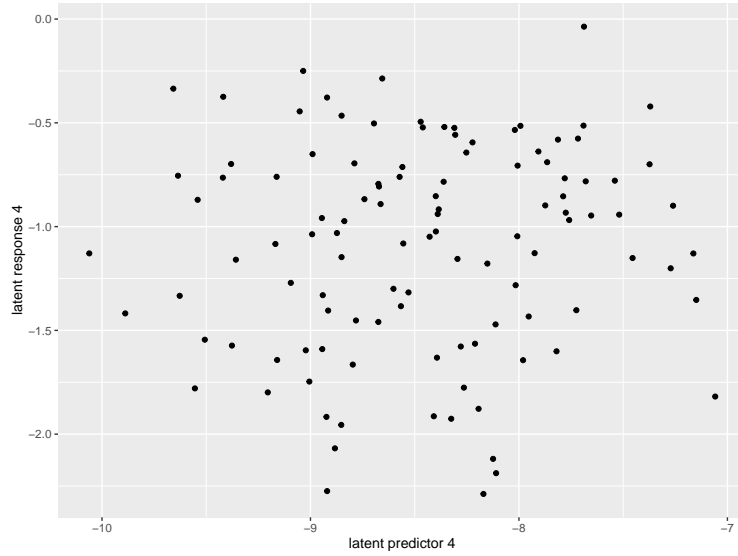


Figure 6: Scatter plots of the 4th latent responses versus the latent predictors for the yeast data estimated by the SOFAR method

### 3. Summary

### References

Rachel B Brem and Leonid Kruglyak. The landscape of genetic complexity across 5,700 gene expression traits in yeast. *Proceedings of the National Academy of Sciences*, 102(5):

1572–1577, 2005.

Ran Dai and Rina Foygel Barber. The knockoff filter for FDR control in group-sparse and multitask regression. *arXiv preprint arXiv:1602.03589*, 2016.

Michael C Gustin, Jacobus Albertyn, Matthew Alexander, and Kenneth Davenport. MAP Kinase Pathways in the Yeast *Saccharomyces cerevisiae*. *Microbiol. Mol. Biol. Rev.*, 62(4):1264–1300, 1998.

Minoru Kanehisa, Susumu Goto, Yoko Sato, Masayuki Kawashima, Miho Furumichi, and Mao Tanabe. Data, information, knowledge and principle: back to metabolism in KEGG. *Nucleic acids research*, 42(D1):D199—D205, 2013.

Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.

Yoshimasa Uematsu, Yingying Fan, Kun Chen, Jinchi Lv, and Wei Lin. SOFAR: large-scale association network learning. *IEEE Transactions on Information Theory*, 2019.

Jianxin Yin and Hongzhe Li. A sparse conditional gaussian graphical model for analysis of genetical genomics data. *The annals of applied statistics*, 5(4):2630, 2011.