

# Principal Components Analysis Based on Unsupervised Feature Extraction Applied to Gene Expression Analysis of Blood from Dengue Haemorrhagic Fever Patients

Y-h. Taguchi<sup>1,\*</sup>

<sup>1</sup>Department of Physics, Chuo University, Tokyo, 112-8551, Japan

\*tag@granular.com

## ABSTRACT

### 1 Theoretical background of PCA based unsupervised FE

Although it was empirically established that PCA based unsupervised FE worked well for a wide range of FEs/FSs when applied to gene expression/epigenetic profiles<sup>1-14</sup>, the lack of theoretical background or justification has prevented other researchers from employing this methodology widely. It also prevented us from estimating in which circumstances it works well *a priori* (i.e., before applying this methodology to the specific problem). Here, we propose the theoretical background of this methodology for the first time based upon Ref.<sup>15</sup>, which proved the equivalence between PCA and K-means, although Ding and He<sup>15</sup> did not recognize that their theoretical framework can be applicable to FS, because they applied PCA only to embedded samples, not to embedded features. In their paper, they proposed the  $K$  non-negative indicator vector  $H_K = (\mathbf{h}_1, \dots, \mathbf{h}_K)$ , where

$$\mathbf{h}_k = (0, \dots, 0, \overbrace{1, \dots, 1}^{n_k}, 0, \dots, 0)^T / n_k^{1/2}$$

and  $n_k$  is the size of the  $k$ th cluster, the elements with  $h_k = 1$  belong to the  $k$ th cluster. They also showed that the connectivity matrix  $C$  is defined and represented as

$$C = \sum_{k=1}^K \mathbf{h}_k \mathbf{h}_k^T = \frac{1}{N} \mathbf{e} \mathbf{e}^T + \sum_{k=1}^{K-1} \mathbf{u}_k \mathbf{u}_k^T,$$

where  $\mathbf{e} = (1, \dots, 1)^T$ . Thus, using PC scores, we could derive cluster structures among genes. In other words, embedding features by PCA is equivalent to figuring out how genes are clustered in the fully unsupervised manner. By computing  $C$  of synthetic data ( $s = 2$ ) with  $K = 1$ , we could identify that the  $C$  between genes  $i = 991, \dots, 1000$ , i.e., those with gene expression distinct between the two classes, have larger  $C$  values (see Fig. S6). This suggested that the theory proposed by Ref.<sup>15</sup> could be applicable not only to sample classification, as they have done, but also to FSs, as has been demonstrated in PCA based unsupervised FE. Thus, Ding and He<sup>15</sup> provided the theoretical background as to why PCA based unsupervised FE works well; PCA based unsupervised FE tries to cluster genes that share similar expression profiles and this results in the selection of genes with expressions distinct between two classes. One may wonder why we do not use directly K-means instead of PCA, if they are equivalent. The reason is that there is only one small cluster to which a limited number of (in this case, as small as 10) genes belong, while the majority (99 % genes) do not form any clusters. K-means must cluster all genes without exceptions. This forces K-means to generate non-existent clusters. Actually, if we apply K-means assuming two clusters, we could never obtain a small cluster including only 10 genes, instead we have two broad clusters whose sizes are equivalent to each other. This kind of methodological limitation exists in all clustering methodologies, because all clustering methodology must cluster elements into a limited number of clusters, even when there is only one small cluster to which a very small part of the elements belong and majority do not form any clusters at all. To the best of our knowledge, PCA based unsupervised FE is the only methodology that could deal with this kind of difficult-to-treat situation. This is possibly the reason why PCA based unsupervised FE could outperform other conventional methodologies for a wide range of problems. In addition, Ding and He<sup>15</sup> provided the missing criteria concerning in which circumstances PCA based unsupervised FE is recommended. Simply speaking, if there is a limited number of small clusters to which a limited proportion of elements belong while the majority of elements do not form any

clusters, PCA based unsupervised FE is useful. To see if this is the case for a DENV data set, we computed  $C$  for data set 2 with  $K = 3$  and ordered columns/rows using the spectral ordering<sup>16</sup> (Fig. S6). The results definitely showed that situation is even worse; there are no clear clusters (block diagonal parts) although around the corners there are some genes with high connectivities, because if there are clusters, multiple block diagonal structures should appear as demonstrated by Refs.<sup>15,16</sup>. This is possibly why PCA based unsupervised FE could outperform those methodologies that cannot deal with this situation effectively. It also supports the employment of two minor PCs (2nd and 3rd) that result in the clear appearance of a set of genes associated with high connectivity. Furthermore, we computed the correlation coefficient between  $\mathbf{q}_1$ , that is the *continuous* inverse index permutation<sup>16</sup> and is used to order features, and  $P$ -values attributed to each gene by PCA based unsupervised FE. Correspondingly, high correlation (Pearson's correlation coefficient is as large as 0.627, that is associated with  $P < 2.2 \times 10^{-16}$ ) was observed. This supported the use of  $P$ -values for FS instead of  $\mathbf{q}_1$ , which requires diagonalization of an  $N \times N$  (thus, generally huge) matrix. We believe the discussion in this subsection justifies the use of PCA based unsupervised FE for the difficult situation where there are only a few (or even no) clusters to which a limited number of elements belong while the majority of elements do not form any clusters, which was a difficult situation that could not be dealt with well by other methods. For more details about the computation of  $C$ , see below.

## 2 Additional methodological advantages of PCA based unsupervised FE

In the previous subsection, we discussed the general methodological advantages of PCA based unsupervised FE based upon the theory proposed by<sup>15</sup>. There are several additional advantages of PCA based unsupervised FE. For example, one may wonder why genes were not screened directly based on the criteria used to specify the PCs for FEs, i.e., DHF+DF vs. CP+HC or convalescent vs. acute. Other than the problem that there are too many genes identified (see above), selecting genes because of their fitness to assumed categorical classes is problematic. It was impossible to reconstruct a two-dimensional space where DHF and DF were well discriminated, as has been done in the present research. Fig. S7 shows the distribution of genes attributed to the two classes on the plane spanned by the PC2 and PC3 loadings. It is obvious that the genes are not unidirectionally distributed around the origin, but alongside the diagonal directions; this is the direction that mostly represents the distinction between the two classes. This means that to construct a two-dimensional space where DF and DHF were well discriminated, we need to include genes unrelated to the distinction between the two classes. This shows the limitation of the supervised method, which can select something targeted, while unsupervised FE can depict something not intended but related to the critical biological background. No supervised method can overcome this difficulty because they cannot select genes that are not specific to something targeted. It is unrealistic to assume that we know everything; therefore, a supervised method might miss something biologically important unintentionally. Thus, unsupervised methods are preferable to supervised if the unsupervised method can be applied to the data set.

Another advantage of the unsupervised method is the number of classes that should be assumed when FE is performed. Although data sets 1 and 2 apparently comprise four classes, in our analysis, we identified that two classes is a reasonable assumption. However, it is difficult for supervised methods to assume a suitable number of classes, because the number of classes is not supposed to be identified, but to be assumed by supervised methodology. Thus, it is evident that assuming two classes not four classes in data sets 1 and 2 is the reason of the successful FEs, and unsupervised FE is more suitable to the present study than supervised FEs.

Furthermore, although PCA is supposed not to be able to represent non-linearity, because PCA is a linear method, this is not always true. For example, in Fig. 4 in main text, development time of diseases is not proportional to any of the gene expressions, because it curves. However, since PCA identifies a two-dimensional space where non-linearity can be expressed as a curve, PCA identified successfully the non-linear dependence of development time upon gene expression. In this sense, if PCA could detect more than one-dimensional space, e.g., a plane, non-linearity could be captured, even using linear methods like PCA.

## 3 Details about sam and limma

When using sam, gene expression is given to `sam` assuming two or four class arrangements. Then probes associated with `q.value` less than 0.01 were identified as selected genes. When using limma assuming two classes, pseudo R code is

```
gene_exp <- new("ExpressionSet", expr=data.matrix(log(x[, -1])))
fData(gene_exp)[["gene_id"]] <- x[, 1]
pData(gene_exp)[["sample_name"]] <- class
design <- model.matrix(~0+class)
colnames(design) <- levels(class)
fit <- lmFit(gene_exp, design)
fit <- eBayes(fit)
```

```

#----- four two classes -----
TT <- topTableF(fit, adjust="BH", number=dim(x)[1])
#----- for four classes -----
contrast.matrix <- makeContrasts(class1-class2, class1-class3, class1-class4, class2-class3,
  class2-class4, class3-class4, levels=design)
fit2 <- contrasts.fit(fit, contrast.matrix)
fit2 <- eBayes(fit2)
TT <- topTableF(fit2, adjust="BH", number=dim(x)[1])
#-----
table(TT[, 6] < 0.01)

```

where  $x$  is supposed to include gene expression, with rows and columns being genes and samples, respectively. The first column is supposed to be gene identifier. `class` is supposed to be factor that represents sample classes. `TT[, 6]` is supposed to include adjusted  $P$ -values.

## 4 Details of computing connectivity matrix $C$

### 4.1 Synthetic data set

For synthetic data set, a  $1000 \times 1000$  matrix was generated as described in main text. Then,  $C$  is computed. Fig. S6 shows the connectivity matrix between 900th and 1000th genes. Only genes between 990th and 1000th are associated with distinct expression between two classes.

### 4.2 Data set 2

After computing connectivity matrix  $C$ , the eigen vector  $\mathbf{q}_1$  was computed. Starting initial random vector  $\mathbf{q}_1$  drawn from the uniform distribution  $(0, 1]$ , only three iterations of  $\mathbf{q}_1 \leftarrow C\mathbf{q}_1$  with suitable scaling  $|\mathbf{q}_1| = 1$  turned out to be enough for the convergence. Since  $-\mathbf{q}_1$  is also an eigen vector if  $\mathbf{q}_1$  is an eigen vector, we could not identify if gene are ordered in the decreasing or increasing order of the elements of  $\mathbf{q}_1$ , we first compute the correlation between  $\mathbf{q}_1$  and  $P$ -values that PCA based unsupervised FE attributed to each gene. Then, genes are ordered such that those associated with smaller  $P$ -values are top ranked. Then, connectivities among top ranked 2400 genes are drawn in Fig. S6 after averaging over every 10 sequentially ranked genes.

## References

1. Taguchi, Y.-h. Identification of aberrant gene expression associated with aberrant promoter methylation in primordial germ cells between E13 and E16 rat F3 generation vinclozolin lineage. *BMC Bioinformatics* **16**, S16 (2015). URL <http://www.biomedcentral.com/1471-2105/16/S18/S16>.
2. Taguchi, Y.-h. Integrative analysis of gene expression and promoter methylation during reprogramming of a non-small-cell lung cancer cell line using principal component analysis-based unsupervised feature extraction. In Huang, D.-S., Han, K. & Gromiha, M. (eds.) *Intelligent Computing in Bioinformatics*, vol. 8590 of *LNCS*, 445–455 (Springer International Publishing, Heidelberg, 2014).
3. Taguchi, Y.-h., Iwadate, M., Umeyama, H., Murakami, Y. & Okamoto, A. Heuristic principal component analysis-based unsupervised feature extraction and its application to bioinformatics. In Wang, B., Li, R. & Perrizo, W. (eds.) *Big Data Analytics in Bioinformatics and Healthcare*, 138–162 (2015).
4. Taguchi, Y.-H., Iwadate, M. & Umeyama, H. Heuristic principal component analysis-based unsupervised feature extraction and its application to gene expression analysis of amyotrophic lateral sclerosis data sets. In *Computational Intelligence in Bioinformatics and Computational Biology (CIBCB), 2015 IEEE Conference on*, 1–10 (2015).
5. Taguchi, Y. H., Iwadate, M. & Umeyama, H. Principal component analysis-based unsupervised feature extraction applied to in silico drug discovery for posttraumatic stress disorder-mediated heart disease. *BMC Bioinformatics* **16**, 139 (2015).
6. Umeyama, H., Iwadate, M. & Taguchi, Y. H. TINAGL1 and B3GALNT1 are potential therapy target genes to suppress metastasis in non-small cell lung cancer. *BMC Genomics* **15 Suppl 9**, S2 (2014).
7. Murakami, Y. *et al.* Comprehensive analysis of transcriptome and metabolome analysis in Intrahepatic Cholangiocarcinoma and Hepatocellular Carcinoma. *Sci Rep* **5**, 16294 (2015).
8. Murakami, Y. *et al.* Comparison of Hepatocellular Carcinoma miRNA Expression Profiling as Evaluated by Next Generation Sequencing and Microarray. *PLoS ONE* **9**, e106314 (2014).

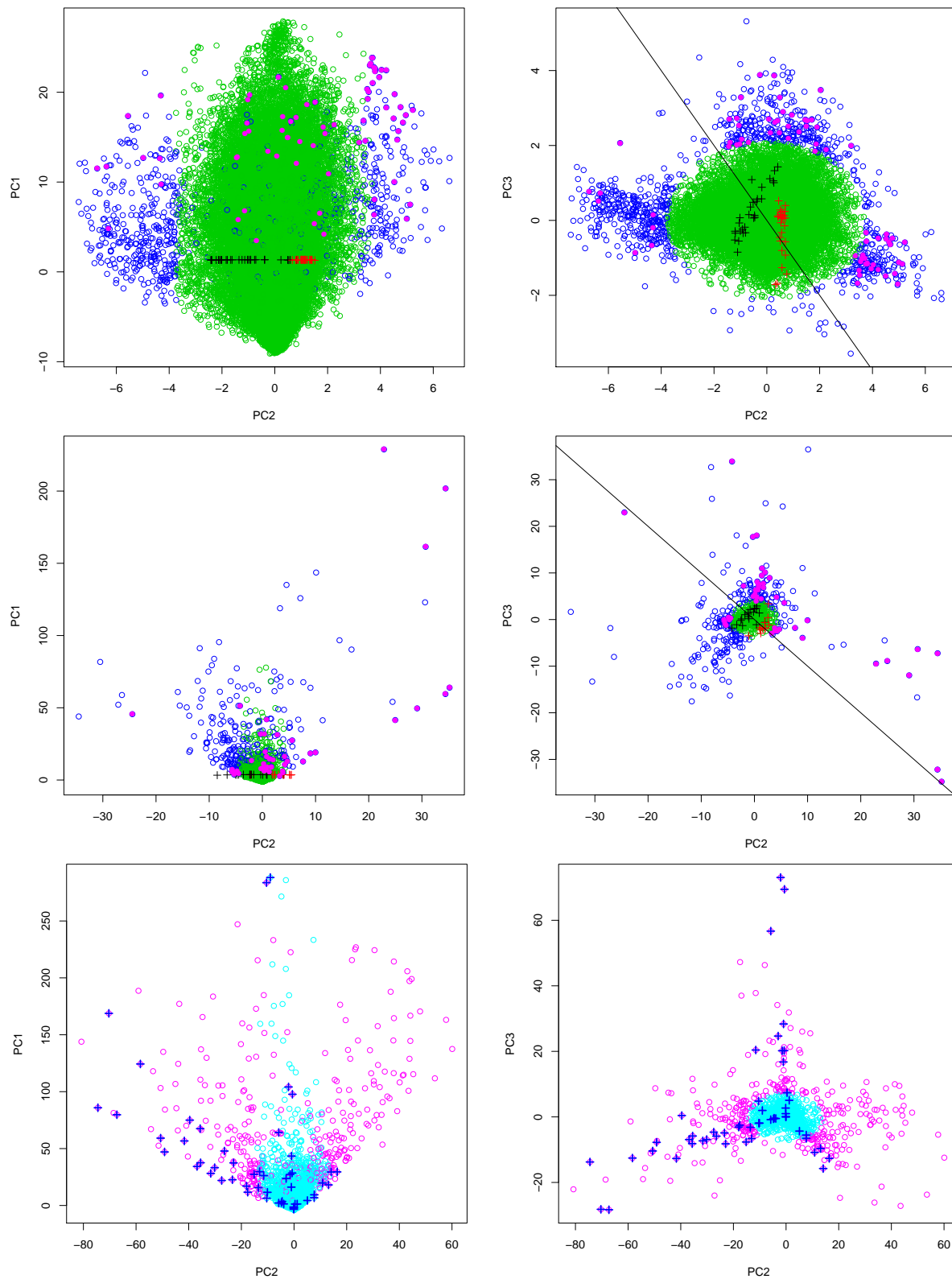
9. Murakami, Y. *et al.* Comprehensive miRNA expression analysis in peripheral blood can diagnose liver disease. *PLoS ONE* **7**, e48366 (2012).
10. Taguchi, Y. H. & Murakami, Y. Universal disease biomarker: can a fixed set of blood microRNAs diagnose multiple diseases? *BMC Res Notes* **7**, 581 (2014).
11. Taguchi, Y. H. & Murakami, Y. Principal component analysis based feature extraction approach to identify circulating microRNA biomarkers. *PLoS ONE* **8**, e66714 (2013).
12. Kinoshita, R., Iwadate, M., Umeyama, H. & Taguchi, Y. H. Genes associated with genotype-specific DNA methylation in squamous cell carcinoma as candidate drug targets. *BMC Syst Biol* **8 Suppl 1**, S4 (2014).
13. Ishida, S., Umeyama, H., Iwadate, M. & Taguchi, Y. H. Bioinformatic Screening of Autoimmune Disease Genes and Protein Structure Prediction with FAMS for Drug Discovery. *Protein Pept. Lett.* **21**, 828–39 (2014).
14. Taguchi, Y.-h. & Okamoto, A. Principal component analysis for bacterial proteomic analysis. In Shibuya, T., Kashima, H., Sese, J. & Ahmad, S. (eds.) *Pattern Recognition in Bioinformatics*, vol. 7632 of *LNCS*, 141–152 (Springer International Publishing, Heidelberg, 2012).
15. Ding, C. & He, X. K-means clustering via principal component analysis. In *Proceedings of the Twenty-first International Conference on Machine Learning, ICML '04*, 29– (ACM, New York, NY, USA, 2004). URL <http://doi.acm.org/10.1145/1015330.1015408>.
16. Ding, C. & He, X. Linearized cluster assignment via spectral ordering. In *Proceedings of the Twenty-first International Conference on Machine Learning, ICML '04*, 30– (ACM, New York, NY, USA, 2004). URL <http://doi.acm.org/10.1145/1015330.1015407>. DOI:10.1145/1015330.1015407.

**Table S1.** List of samples included in data set 1, 2 and 3. DSS:Dengue Shock Syndrome. GSE51808: RMA normalization was performed using Expression Console software. GSE13052: Intensity was acquired using Beadstudio software Intensity was background normalised (Subtract the background value). GSE25001: Data was normalised by Beadstudio software. GSE9378: Signal values were calculated using robust multi-array analysis (RMA) (BioConductor), transformed using inverse nlog, and then imported into GeneSpring (Agilent) for chip normaliaztion to 50th percentile and gene normalization to the mean of controls (where available) for each cell type independently. GSE43777: RMA normalization was performed using Expression Console software. For more details, see paper.

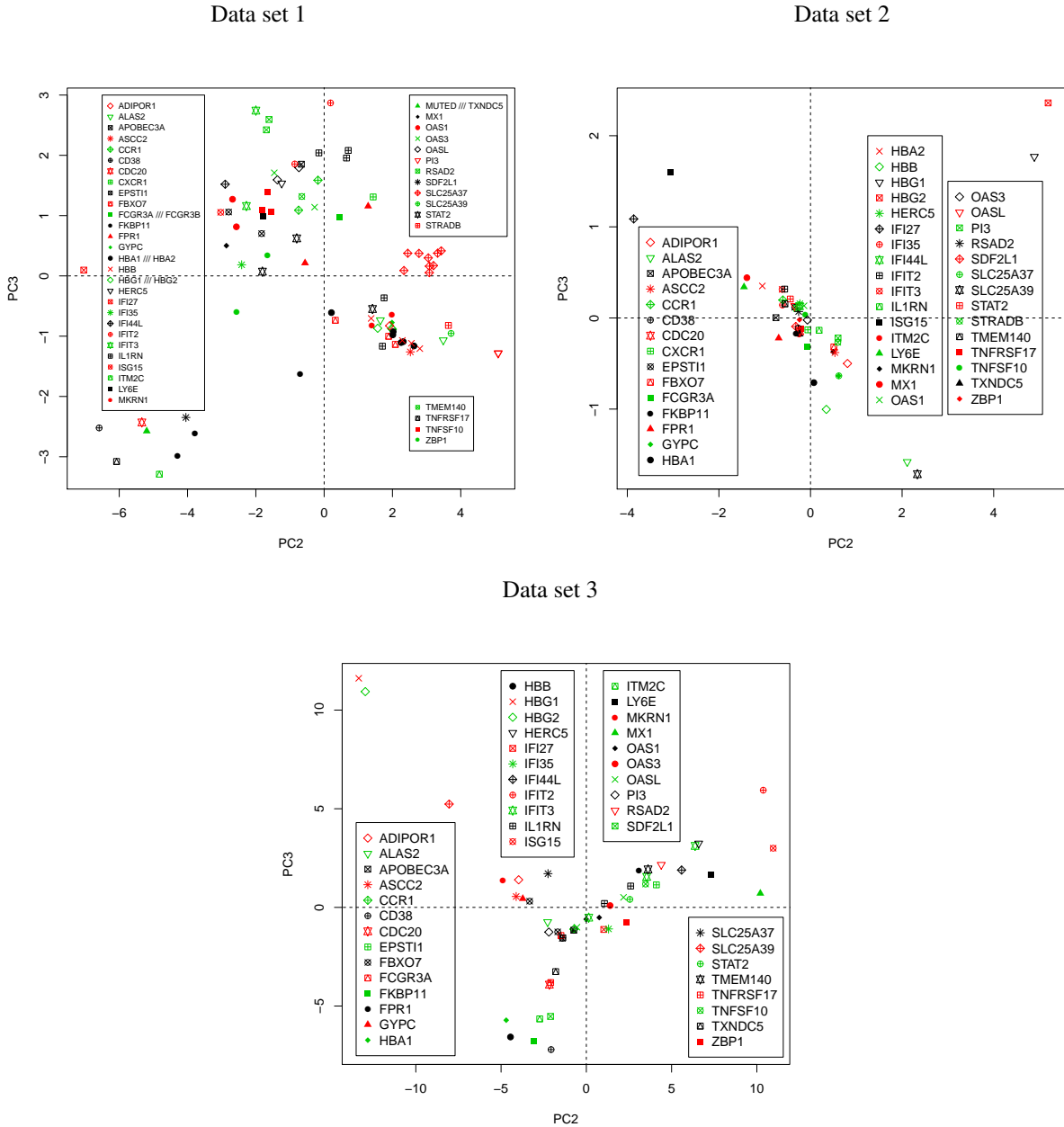
Data set 1 (GSE51808)		Affymetrix HT HG-U133+ PM Array Plate							
	Healthy Controls (HC)	Acute Patients (AC)	DF	DHF					
	9	19	18	10					
Data set 2 (GSE13052)		Sentrix HumanRef-8 Expression BeadChip							
	Acute	Convalescent							
uncomplicated (DF)	10	5							
DSS* (DHF)	9	6							
Data set 3 (GSE25001)		Illumina humanRef-8 v2.0 expression beadchip							
	Acute	0-1	Disease (Fever)	follow up					
DF	56	32	31	16					
DHF	24	12	20	18					
<i>in vitro</i> (GSE9378)		Affymetrix Human Genome U133A Array							
	HUVEC	Monocyte							
control	2	2							
infected	2	2							
Data set 4 (GSE43777-GPL570)		Affymetrix Human Genome U133 Plus 2.0 Array							
	G0	G1	G2	G3	G4	G5	G6	G7	
DF	0	2	5	8	9	5	11	12	
DHF	0	0	3	8	10	5	11	12	
Data set 5 (GSE43777-GPL201)		Affymetrix Human HG-Focus Target Array							
DF	2	5	21	18	22	22	24	45	
DHF	0	0	0	1	3	1	1	3	

**Table S2.** Number of genes identified by sam, limma and PCA based unsupervised FE. Two classes mean “DHF+DF” vs “CP+HC” for data set 1 (GSE51808) and “Acute” vs “Convalescent” for data set 2 (GSE13052). \*:all probes. The numbers in parentheses are those when the sample numbers are halved. Averaged values over 100 ensembles are presented. Halving was performed within each of four classes. Thus, the ratio between classes was conserved.

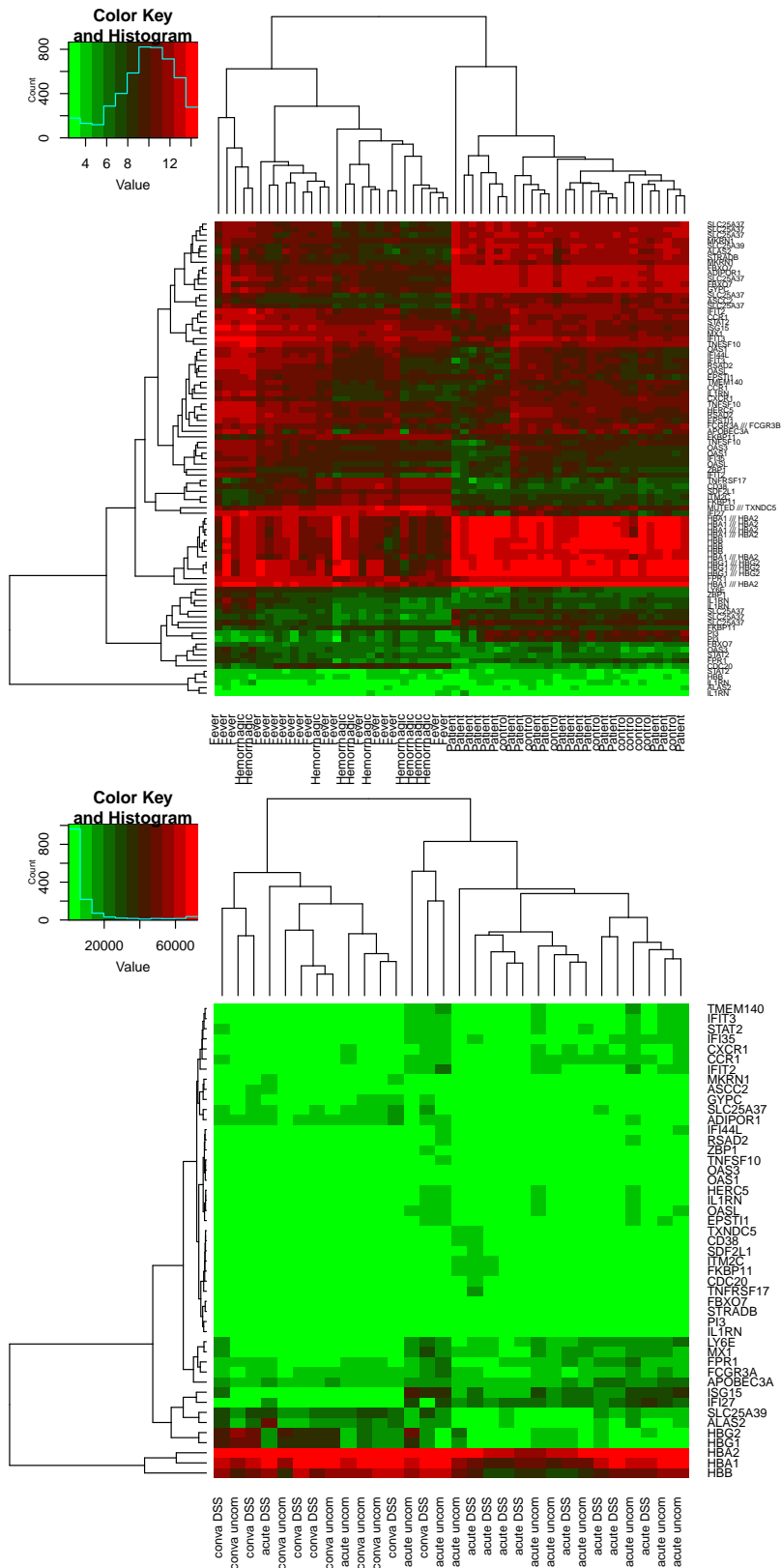
Data set	sam		limma		PCA based unsupervised FE
	two classes	four classes	two classes	four classes	
1	17680 (18469)	16647 (7461)	54715* (54715*)	13506 (5706)	879 (826)
2	2427 (41)	865 (0)	21795 (19855)	20629 (17478)	275 (286)



**Figure S1. Top:** Biplot of PC1 to PC3 for data set 1 (GSE51808). Open green circles are probes not selected as outliers. Open blue circles are 879 probes identified as outliers. Black and red crossed represent patients with symptom (DF/DHF) and those without symptom (HC/AC). Solid line represents the line  $PC2 = -PC3$  that roughly represents the distinction between patients with/without symptom. Open magenta circles are probes associated with 46 genes commonly identified as outliers in data set 1 and 2. **Middle:** Biplot of PC1 to PC3 for data set 2 (GSE13052). Open green circles are probes not selected as outliers. Open blue circles are 275 probes identified as outliers. Black and red crossed represent patients with symptom (acute) and those without symptom (convalescent). Solid line represents the line  $PC2 = -PC3$  that roughly represents the distinction between patients with/without symptom. Open magenta circles are probes associated with 46 genes commonly identified as outliers in data set 1 and 2. **Bottom:** Scatter plot of PC1 to PC3 scores attributed to probes for data set 3 (GSE25001). Open cyan circles are probes not selected as outliers. Open magenta circles are probes identified as outliers. Blue crossed represents probes associated with 46 genes commonly identified as outliers in data set 1 and 2.

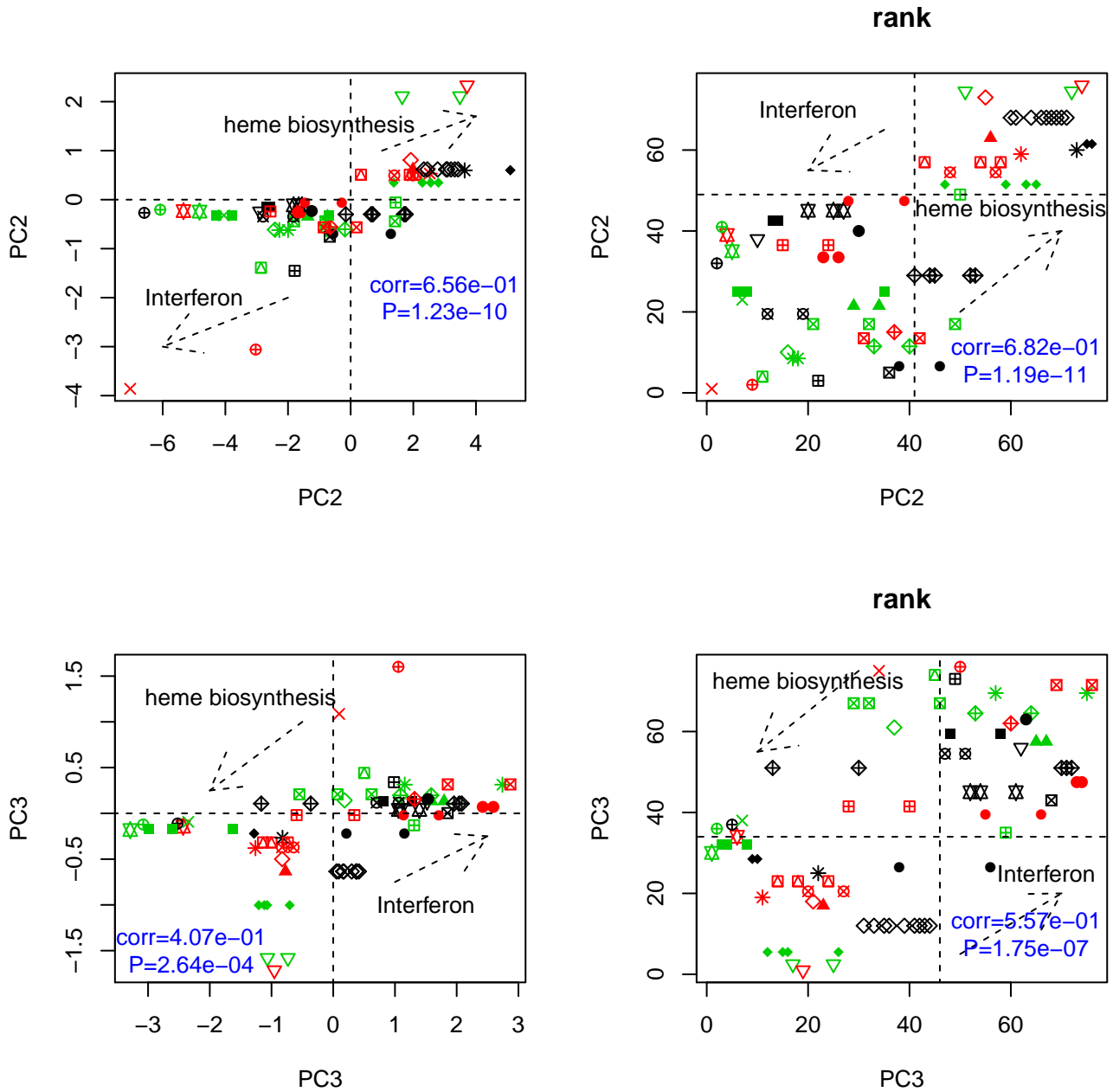


**Figure S2.** Annotation of genes shown in Figs. 4 (upper) and data set 3 in Fig. 5(lower).



**Figure S3.** Heatmaps of probes associated with 46 genes commonly identified in data set 1 and 2. Upper: Data set 1. Samples are strictly grouped as “Fever” (DF) + “hemorrhagic” (DHF) and control (HC) + Patient (CP). Lower: Data set 2. Samples are almost grouped as “conva” (convalescent) and ”acute”. Only two acute and one conva patients were wrongly grouped. Grading: bright red (green) represents more expressive(suppressive) expressions. They are drawn using heatmap function implemented in R. hierarchical clustering was performed by Unweighted Pair Group Method with Arithmetic mean with setting method = ”average” option. Distances between gene expression were Euclidean distance.

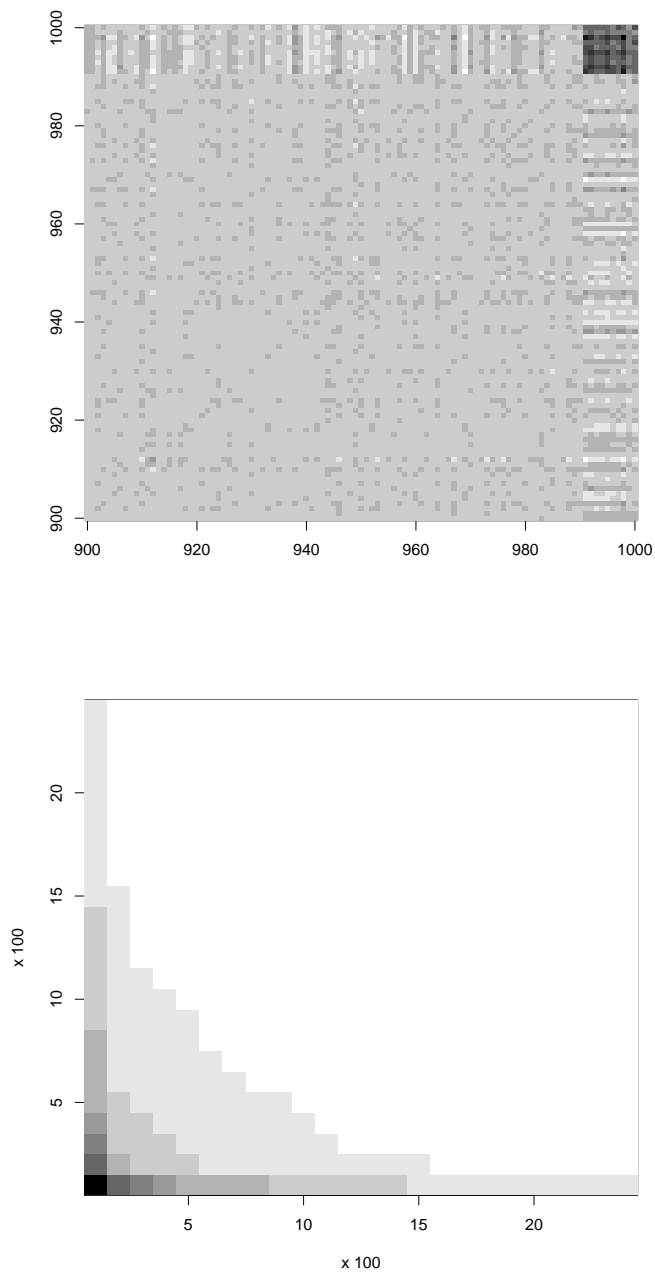




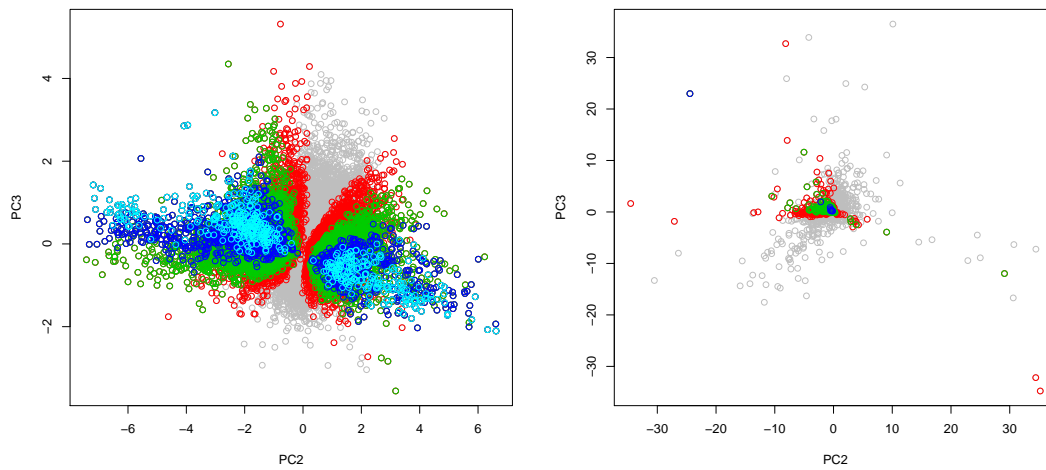
**Figure S4.** Comparison between PC scores attributed to probes associated with 46 genes commonly identified in data set 1 and 2. Upper:PC2, lower:PC3, Left:PC scores, right: rank of PC scores. Correlation coefficients (left:Pearson, right:Spearman), as well as associated  $P$ -values are also shown. For the annotation of characters, see Fig. S5

◇ ADIPOR1	⊕ IL1RN
▽ ALAS2	⊕ ISG15
⊠ APOBEC3A	⊗ ITM2C
* ASCC2	⊠ LY6E
◇ CCR1	⊗ MKRN1
⊕ CD38	⊠ MX1
⊗ CDC20	■ OAS1
⊠ CXCR1	● OAS3
⊗ EPSTI1	▲ OASL
⊠ FBXO7	◆ PI3
■ FKBP11	● RSAD2
● FPR1	× SDF2L1
▲ GYPC	◇ SLC25A37
◆ HBB	▽ SLC25A39
● HERC5	⊠ STAT2
× IFI27	* STRADB
◇ IFI35	◇ TMEM140
▽ IFI44L	⊕ TNFRSF17
⊠ IFIT2	⊗ TNFSF10
* IFIT3	⊠ ZBP1

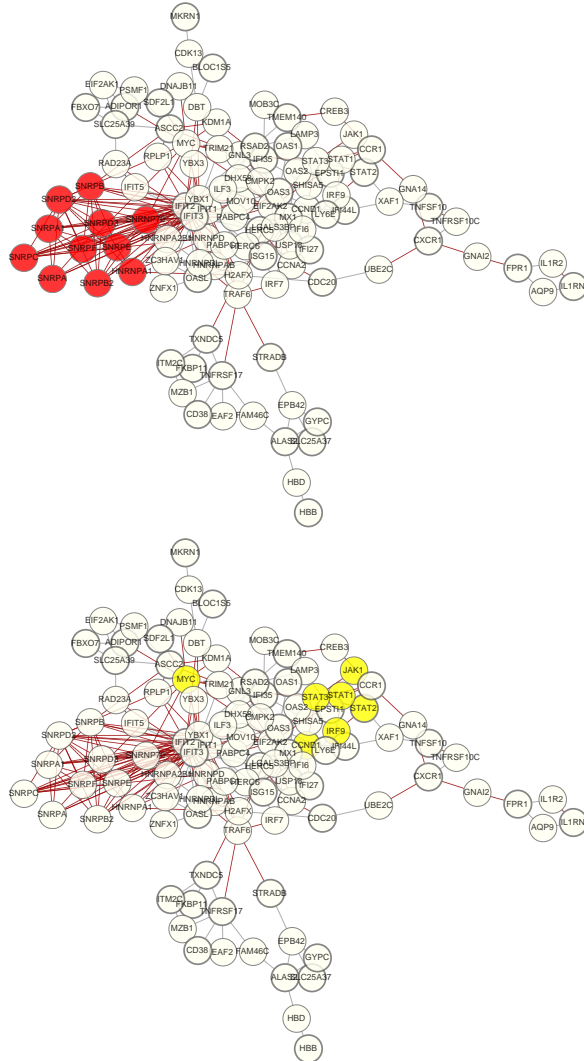
**Figure S5.** Character annotations used in Fig. S4



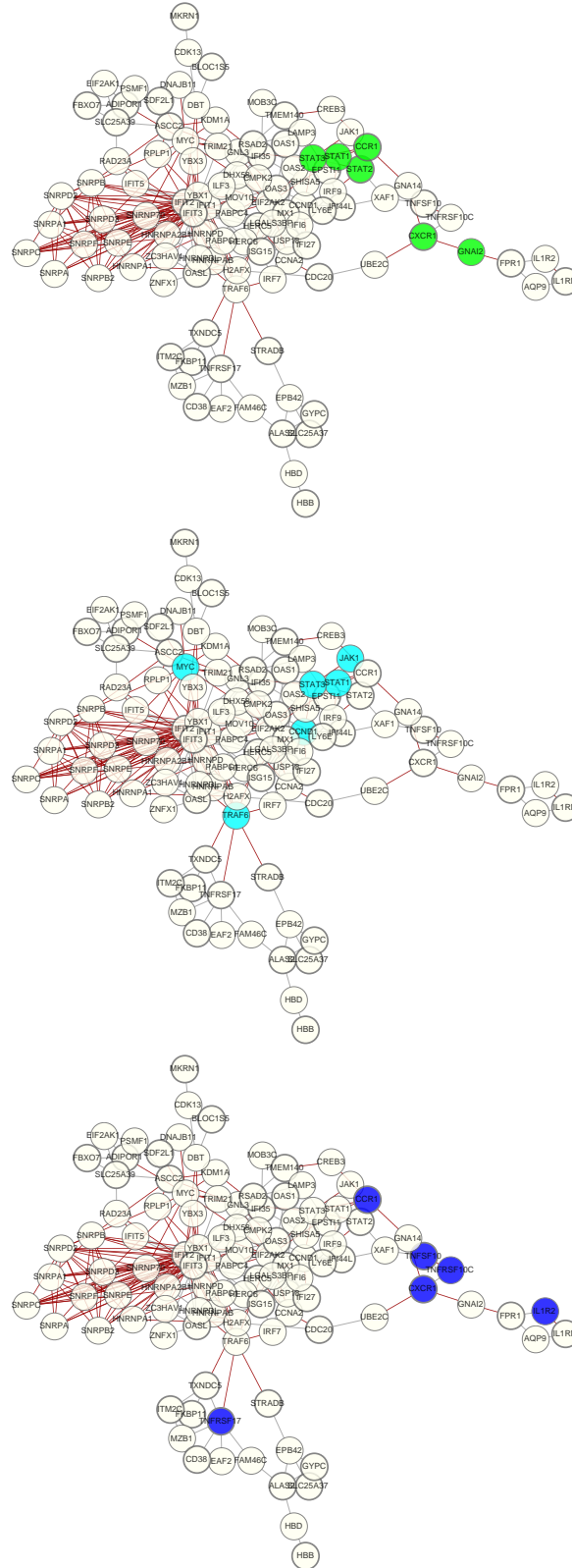
**Figure S6.** Connectivity matrix of synthetic data ( $s = 2$ , left, genes  $900 \leq i \leq 1000$ ) and data set 2 (right, coarse grained values averaged over every 10 sequentially ranked genes among top ranked 2400 genes). Darker gray scales correspond to higher connectivities.



**Figure S7.** Left: Distribution of  $P$ -values adjusted by BH criterion on two dimensional space spanned by PC2 and PC3 loadings for data set 1 (GSE51808).  $P$ -values were computed by  $t$  test and were those to deny null hypothesis that mean of  $x_{ij}$  within DF+DHF are identical to that of  $x_{ij}$  within CP+HC. Red open circles:  $1 \times 10^{-5} < \text{adjusted } P\text{-values} < 0.01$ , Green open circles:  $1 \times 10^{-10} < \text{adjusted } P\text{-values} < 1 \times 10^{-5}$ , blue open circles:  $1 \times 10^{-13} < \text{adjusted } P\text{-values} < 1 \times 10^{-10}$ , cyan open circles:  $\text{adjusted } P\text{-values} < 1 \times 10^{-13}$ . Right: Distribution of same variables but for data set 2 (GSE13052).  $P$ -values were computed by  $t$  test and were those to deny null hypothesis that mean of  $x_{ij}$  within convalescent are identical to that of  $x_{ij}$  within acute. Red open circles:  $1 \times 10^{-3} < \text{adjusted } P\text{-values} < 0.01$ , Green open circles:  $1 \times 10^{-4} < \text{adjusted } P\text{-values} < 1 \times 10^{-3}$ , blue open circles:  $\text{adjusted } P\text{-values} < 1 \times 10^{-4}$ .



**Figure S8.** Co-expression network inferred by COEXPRESSdb. Upper(red):Spliceosome(hsa03040), lower(yellow):Jak-STAT signaling pathway(hsa04630). Genes in bold open circles are 46 genes identified by PCA based unsupervised FE



**Figure S9.** Co-expression network inferred by COEXPRESSdb (continued). Upper(green): Chemokine signaling pathway(hsa04062), middle(cyan): Pathways in cancer(hsa05200), lower(blue) Cytokine-cytokine receptor interaction (hsa04060). Genes in bold open circles are 46 genes identified by PCA based unsupervised FE

