

Group Meeting

2014.09.29

1. Sparse Subspace Clustering
2. Partial Least Squares
3. TCGA Data Analysis

Sparse Subspace Clustering

PCA in a Geometric View:

Given a set of points $\{x_j\}_{j=1}^N$ in R^D , try to find an (affine) subspace $S \subset R^D$ of dimension d , $\dim(S) = d$ that best fits these points.

$$x_j = \mu + U_d y_j + \epsilon_j \quad j = 1, \dots, N$$

where $U_d \in R^{D \times d}$ whose columns form a basis for the subspace and y_j is the vector of new coordinates of x_j in the subspace.

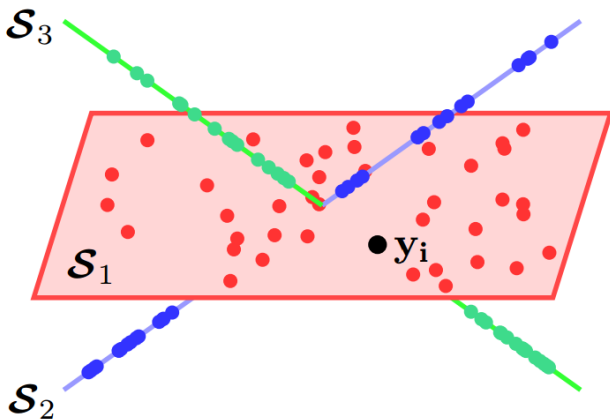
$$\begin{aligned}
& \min_{\mu, U_d, \{y_j\}} && \sum_{j=1}^N \|x_j - \mu - U_d y_j\|^2 \\
& s.t. && U_d^T U_d = I_d, \quad \sum_{j=1}^N y_j = 0
\end{aligned}$$

Subspace Clustering Formulation:

Given a set of sample points $X = \{x_j \in R^D\}_{j=1}^N$ drawn from $n \geq 1$ distinct linear subspaces $S_i \subset R^D$ of dimensions $d_i < D$, $j = 1, \dots, n$.

Identify each subspace S_i without knowing which sample point belong to which subspace.

1. Identifying the number of subspaces n and their dimensions $d_i = \dim(S_i)$
2. Identifying the orthonormal basis for each subspace S_i . (or equivalently a basis for its orthogonal complement S_i^\perp)
3. Clustering the N points into the subspaces to which they belong



How to define a good similarity matrix?

Characterize the local or global **subspace structures** around the points of interest

- Two points could be very close, but lie in different subspaces(e.g. near the intersection of two subspaces)
- Two points could be far from each other, but belong to the same subspace

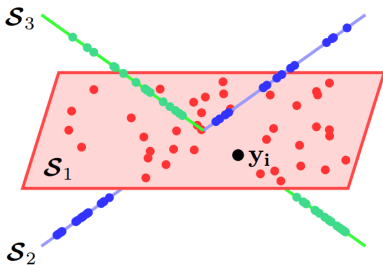
Sparse Subspace Clustering: Key observations and conclusions

Sparse Subspace Clustering: Key observations and conclusions

Self-expressiveness

Each data point in a union of subspaces can be expressed as a (sparse) linear or affine combination of all other points in the dataset.

$$x_j = \sum_{i \neq j} c_{ij} x_i, \quad \text{or} \quad X = XC, \text{diag}(C) = 0$$



Sparse Representation Theory

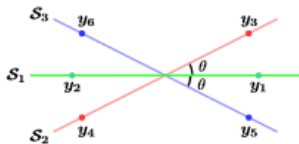
Given the following conditions:

- The subspaces are sufficiently separated
- The data within the subspaces are well distributed

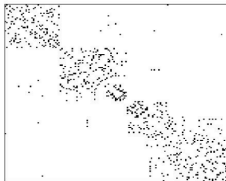
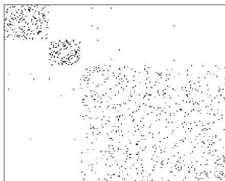
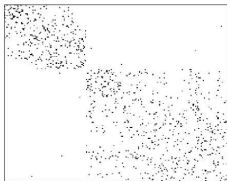
we can recover the sparse subspace representation by the optimization problem:

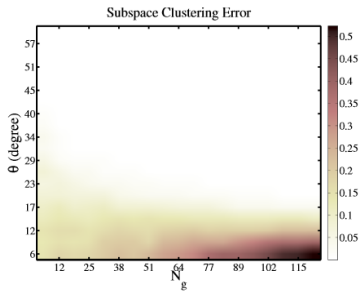
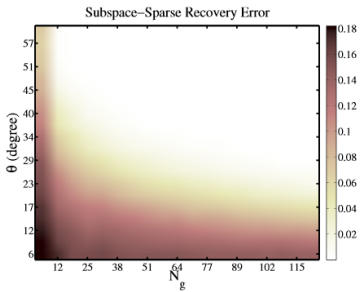
$$\min_C \|C\|_1 \quad \text{s.t. } X = XC, \text{diag}(C) = 0$$

where $\|C\|_1 = \sum_{i,j=1}^N |c_{ij}|$



$$C_1 = \left(\begin{array}{cc|cc|cc} 0 & -1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 & 0 \\ \hline 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & 0 \end{array} \right)$$





Clustering Using Sparse Coefficients

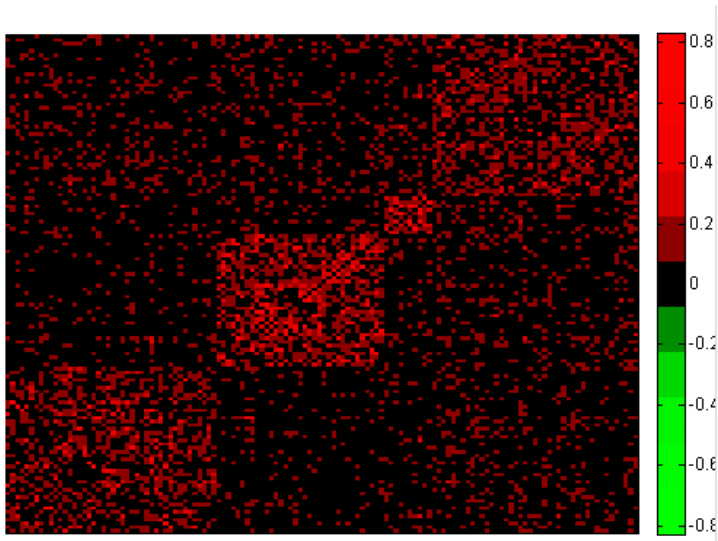
We can use the spectral clustering algorithm.

1. Form a similarity graph $\mathcal{G} = (V, E, W)$ whose nodes are the N data points and whose edges connect points x_i and x_j with a weight $w_{ij} = |c_{ij}| + |c_{ji}|$.

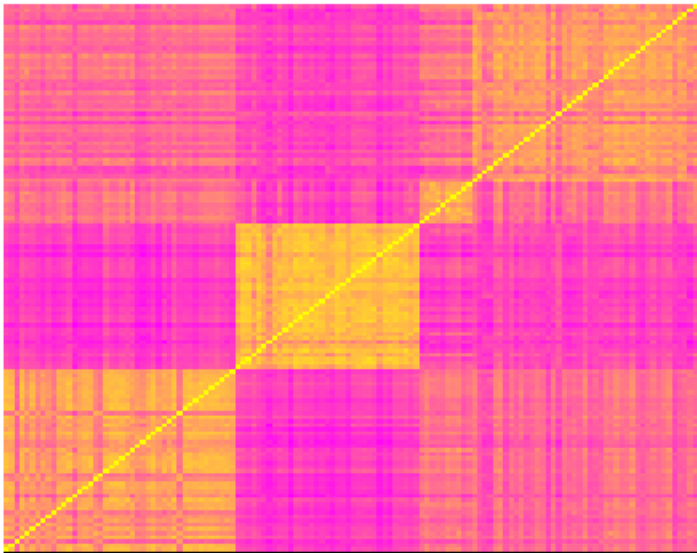
$$W = |C| + |C|^T$$

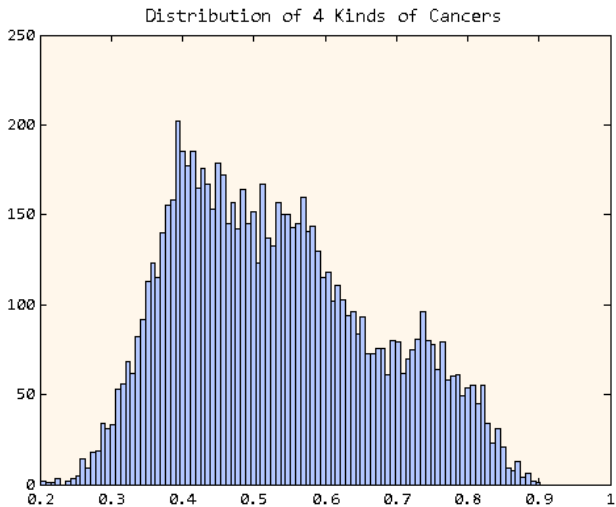
2. Apply spectral clustering algorithm to the similarity graph with weights W to obtain the segmentation of the data.

Application in 4 kinds of Cancers



Correlation Distributions:



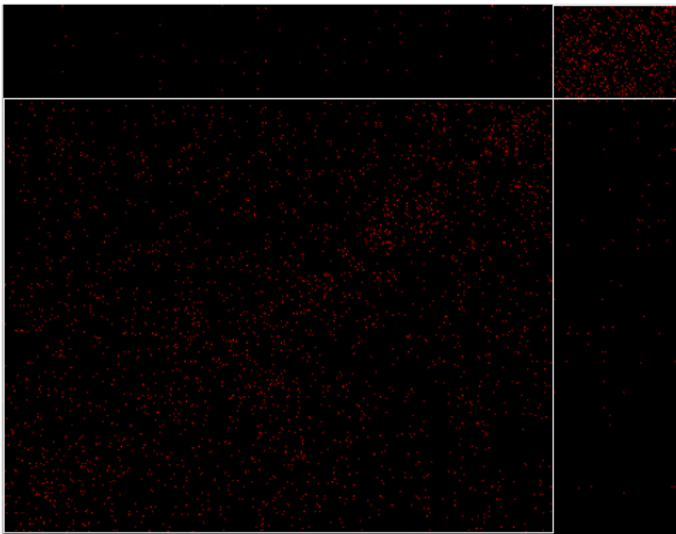


Application in ER+/ER- samples

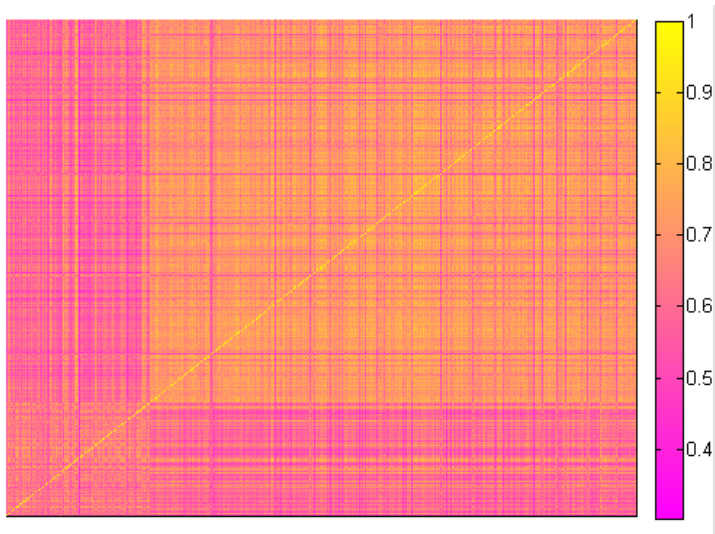
乳腺癌样本数据聚类结果，包含ER+/ER-信息

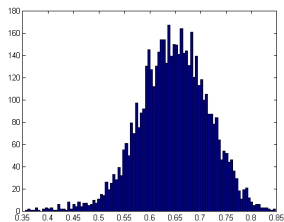
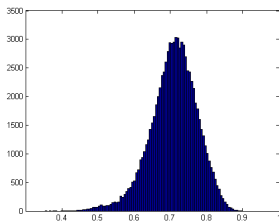
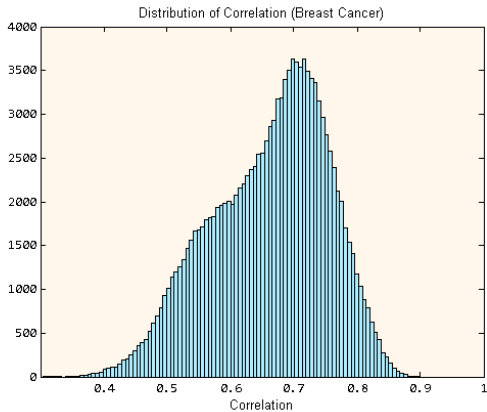
		真实值	
		Positive	Negative
预测值	Positive	370	12
	Negative	24	104

- Elapsed time is 39.406780 seconds



Correlation Heatmap





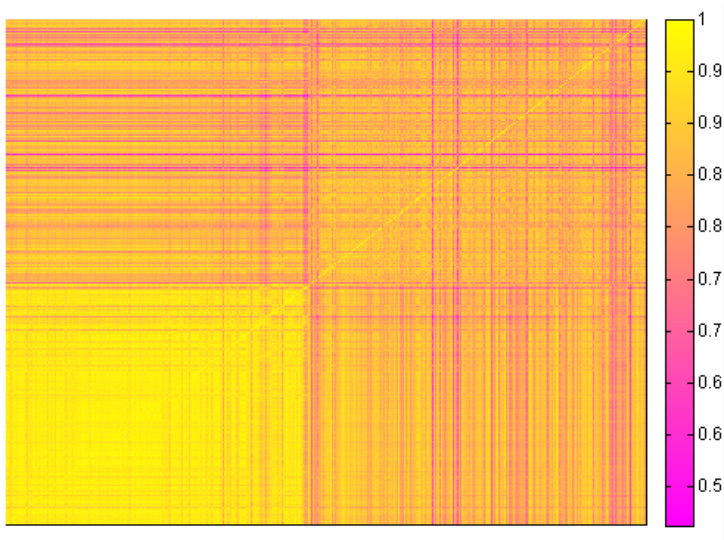
Application in HCC samples

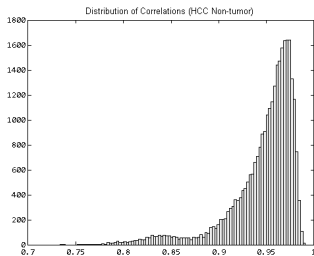
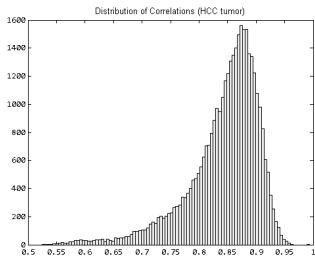
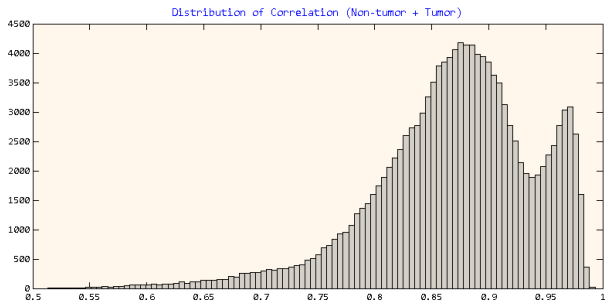
1. GSE25097 Datasets

Label	Healthy	Cirrhotic	Non-tumor	Tumor
Number	6	40	243	268

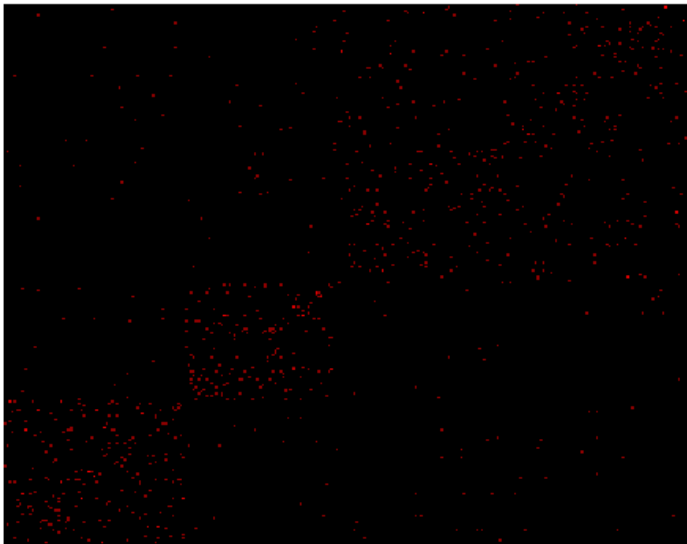


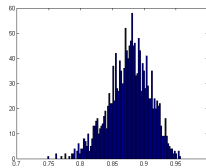
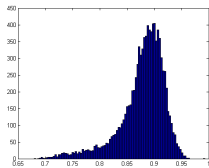
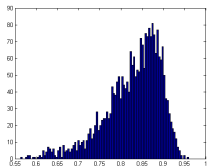
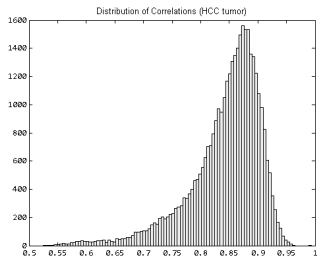
Correlation Heatmap (Non-tumor and tumor)



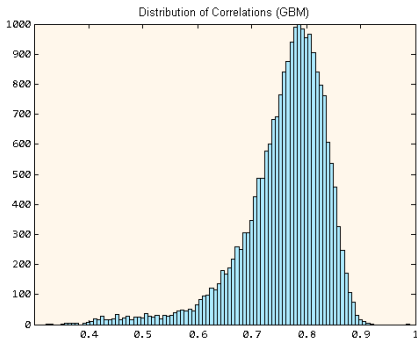


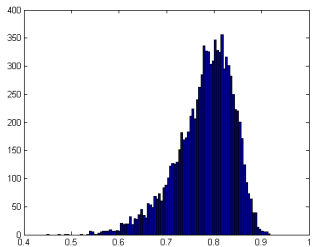
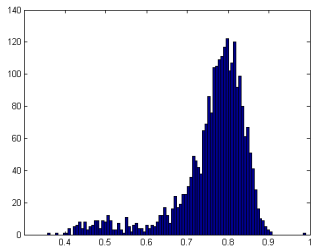
2. Tumor samples only

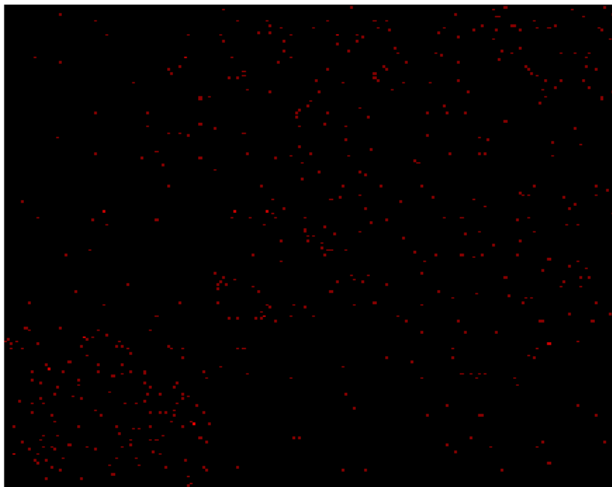




Application in GBM TCGA Dataset







生存期标签 \ 分类结果	1	2
1	33	55
2	38	84
	71	139

Partial Least Squares

PCA in a Statistical View:

Suppose a random variable $x \in R^D$ with zero-mean $\mathbb{E}(x) = 0$, and try to find $d < D$ principle components $y \in R^d$:

- y_i 's are '**uncorrelated**' linear combinations of x :

$$y_i = u_i^T x \in R, u_i \in R^D$$

- The variance of y_i is maximized subject to

$$u_i^T u_i = 1, i = 1, \dots, d \quad \text{Var}(y_1) \geq \text{Var}(y_2) \geq \dots \geq \text{Var}(y_d) > 0$$

Principle Components Regression

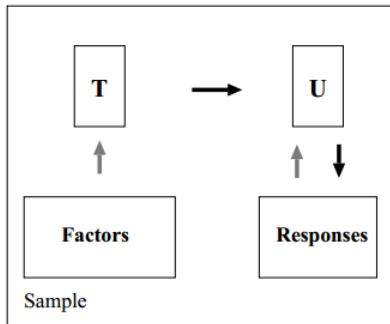
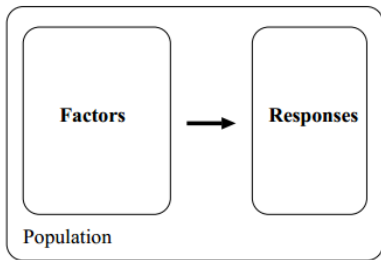
- Form the derived input columns $z_m = Xu_m, m = 1, \dots, M (M \leq p)$
- Regress the response y on z_1, \dots, z_M

Partial Least Squares

Consider both high variance and high correlation with the response:

The m -th PLS direction ψ_m :

$$\begin{array}{ll} \max_{\alpha} & \text{Corr}^2(y, X\alpha) \text{Var}(X\alpha) \text{ or } \text{Cov}(y, X\alpha) \\ \text{s.t.} & \|\alpha\| = 1, \alpha^T \Sigma_N \hat{\psi}_l = 0, l = 1, \dots, m-1 \end{array}$$



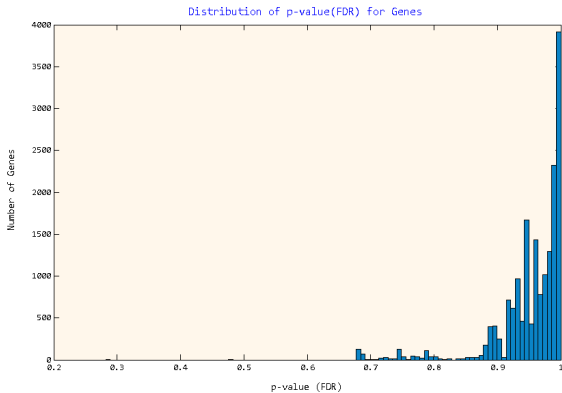
$$X = TP^T + E,$$

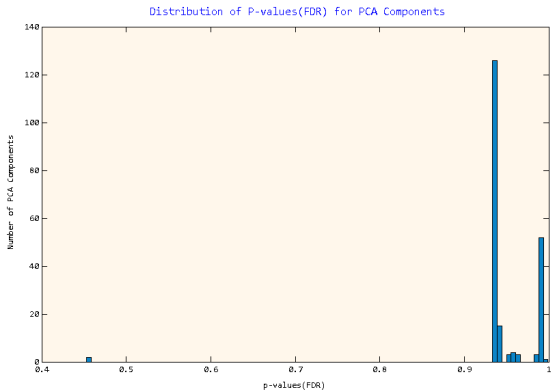
$$Y = UQ^T + F,$$

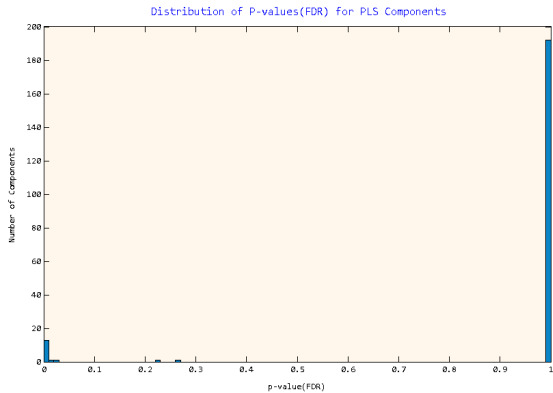
$$U = TB$$

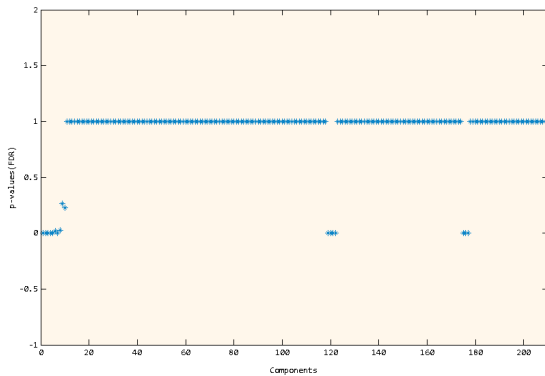
GBM mRNA Data

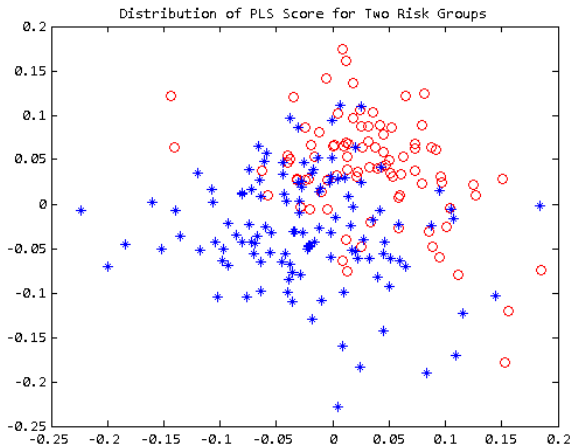
1. Wilcoxon rank sum testing (Survival \sim Gene, PCA, PLS components)

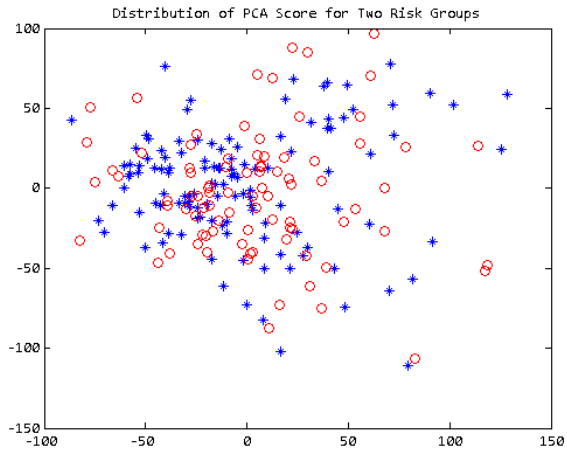




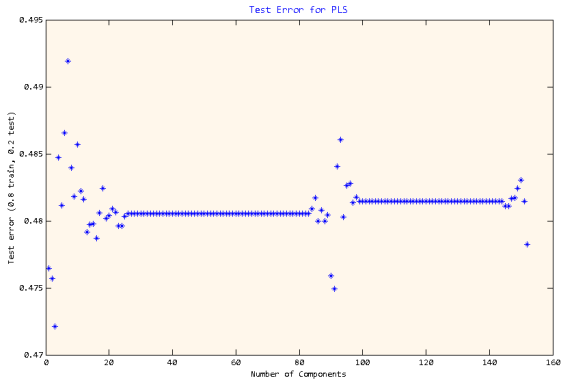








2. Classification error (overfitting)



(Mean error of 100 times)

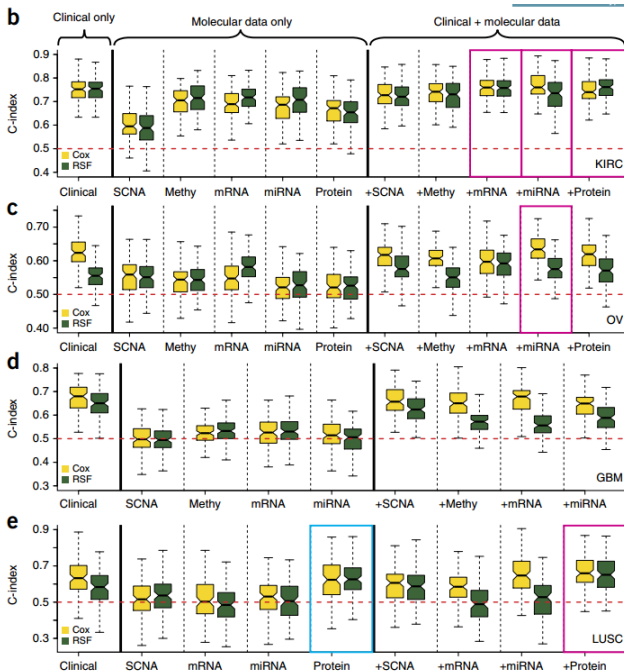
Results of TCGA Survival Data Analysis

how and to what extent TCGA molecular data can affect oncology practice.

ANALYSIS

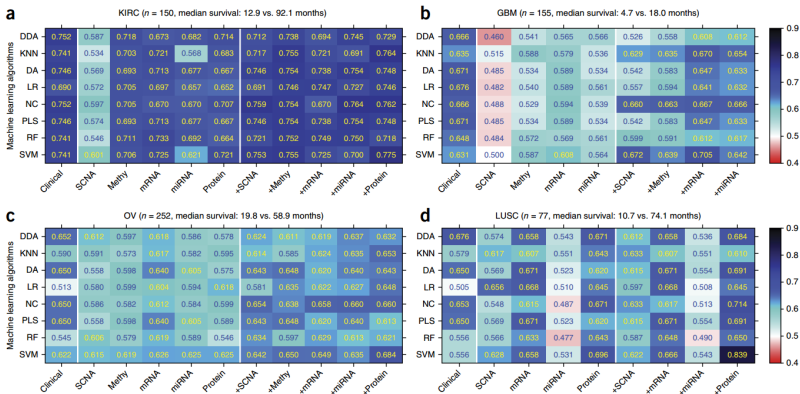
—computational
BIOLOGY

Assessing the clinical utility of cancer genomic and proteomic data across tumor types



- Clinical-variables-only model: show substantial predictive power
- The relative predictive power of individual molecular data sets strongly depended on the cancer type
- Clinical variables + molecular data: significantly improve predictive power but the increase is limited.

Prediction of Dichotomized Survival data



Predictive power of molecular data strongly depends on the cancer type.

