

Clustering of Cancer Data Based on Stiefel Manifold for Multiple Views

Jing Tian

Xinjiang University

Jianping Zhao (✉ zhaojianping@126.com)

Xinjiang University

Chun-hou Zheng

Xinjiang University

Research Article

Keywords: Stiefel manifold, multi-view clustering, cancer data, optimization model, linear search algorithm

Posted Date: February 12th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-154286/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Clustering of Cancer Data Based on Stiefel Manifold for Multiple Views

Jing Tian¹, Jian-ping Zhao^{1,*}, Chun-Hou Zheng^{1,2,1}

¹College of Mathematics and System Sciences, Xinjiang University, Urumqi, China

²School of Computer Science and Technology, Anhui University, Hefei, China

Abstract

Background: In recent years, various sequencing techniques have been used to collect biomedical omics datasets. It is usually possible to obtain multiple types of omics data from a single patient sample. Clustering of these datasets has proved to be valuable for biological and medical research and helpful to reveal data structures from multiple collections. However, such data often have small sample size and high dimension. It is difficult to find a suitable integration method for structural analysis of multiple datasets.

Results: In this paper, a multi-view clustering based on Stiefel manifold method (MCSM) is proposed. Firstly, we established a binary optimization model for the simultaneous clustering problem. Secondly, the optimization problem solved by linear search algorithm based on Stiefel manifold. Finally, we integrated the clustering results obtained from three omics by using k-nearest neighbor method. We applied this approach to four cancer datasets on TCGA. The result shows that our method is superior to several state-of-art methods, which depends on the hypothesis that the underlying omics cluster class is the same.

Conclusion: Particularly, our approach has better performs when the underlying clusters are

¹Corresponding authors

23 inconsistent. For patients with different subtypes, both consistent and differential clusters can be
24 identified at the same time.

25 **Keywords:** Stiefel manifold; multi-view clustering; cancer data; optimization model; linear search
26 algorithm

27 **1. Introduction**

28 One of the challenges of cancer treatment is how to identify tumor subtypes, which can help to
29 provide patients with specific treatment. Meanwhile, with the continuous development of all kinds
30 of sequencing technologies, a lot of high flux data have been produced (Zheng et al. 2016). For
31 cancer subtypes identification, integration of different types of omics data to unravel the molecular
32 mechanism of complex diseases becomes more and more important (Zhang et al., 2016). On the one
33 hand, multiple omics data of different subtypes of cancer provided more detailed information. On
34 the other hand, it led to the new challenge for data analysis. Different levels of multiple omics
35 data often show different types, they have different correlation structure statistical properties and
36 expressions (Zhang et al., 2011). In addition, the same tumor specimens from different levels of data
37 are also unlikely to be independent. Therefore, how to reasonably integrate the multiple omics data
38 to accurately predict cancer subtypes becomes a challenging and interesting research (Bickel et al.,
39 2009).

40 Recently, many strategies have emerged to integrate multi-omics data for association studies. The
41 main objectives of these studies usually are to understand the inter-relationships among various
42 omics data and to detect phenotype-related modules (Manko et al., 2010; Kumar et al., 2011; Yuan
43 et al., 2011). These Multi-view data integration methods can be divided into three categories: early

44 integration, late integration and intermediate integration. Early integration mainly merged all data
45 types into a single dataset and analyzed them directly (Manko et al., 2010). Late integration is the
46 process of analyzing each data type, then combining the results for integration (Fridly et al., 2012).
47 Intermediate integration converts each data type into an intermediate representation, which is then
48 consolidated for analysis (Cancer Genome Atlas Research Network et al., 2012).

49 In recent years, many new methods are also published for clustering. For example, Similarity
50 Network Fusion (SNF) method (Wang et al., 2014) constructed the similarity network for each data
51 type, and then used the iterative method to fuse them into a similar network. The final clustering is
52 obtained by spectral clustering of fusion networks. Some multi-view clustering methods based on
53 spectral clustering have also been proposed (Huang et al., 2012; Kumar et al., 2011; Zhang et al.,
54 2015). They used different integration methods to combine the spectral clustering results from a
55 single view. The Affinity Aggregation for Spectral Clustering (AASC) algorithm (Huang et al., 2012)
56 introduced weights in the spectral clustering of each view, and then added them together to optimize
57 the weights in the calculation.

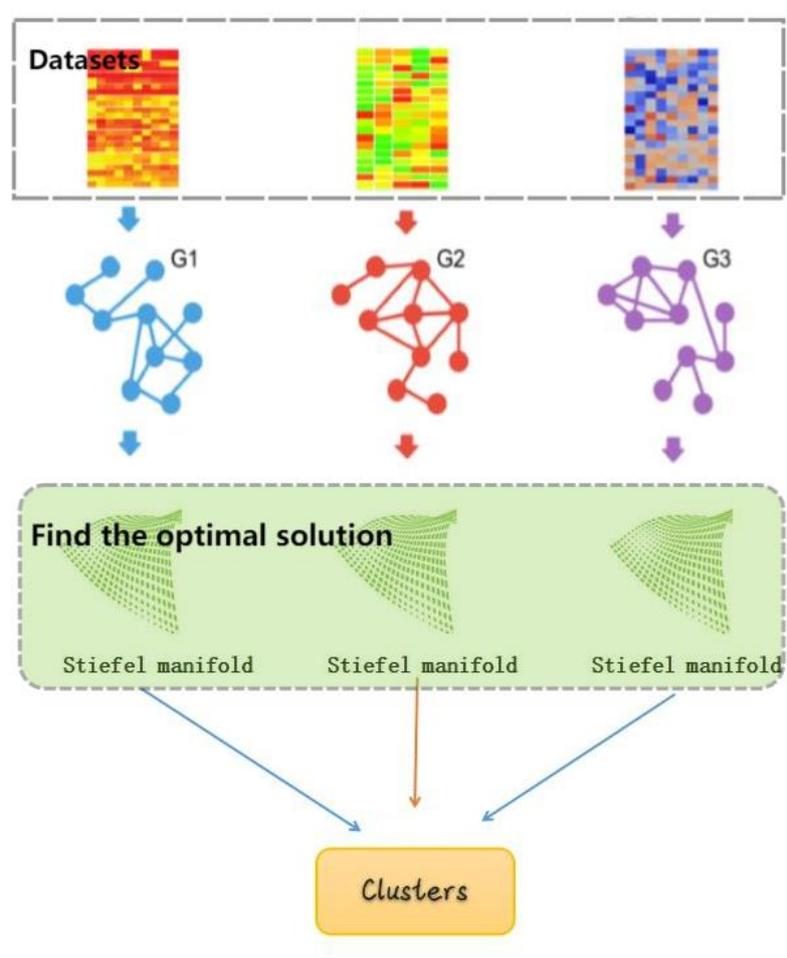
58 However, these methods were put forward based on a basic hypothesis that the underlying omics
59 cluster is the same. In actual situation, there are differences in clusters except for the same cluster
60 (Yuan et al., 2011). In the process of integrated clustering, data clustering was carried out for each
61 view and cluster alignment was carried out for different views, which could handle this situation
62 (Kumar et al., 2011; Zhang et al., 2015; Chen et al., 2017). However, the method (Kumar et al.,
63 2011) tended to obtain the local optimal as described above, and the methods (Zhang et al., 2015;
64 Chen et al., 2017) relax excessively the original multi-view point specific tangent condition, so that
65 the information of each viewpoint may be lost. In the paper (Yu et al., 2019), the authors proposed

66 the Multi-View Clustering using Manifold Optimization (MVCMO) method considering the
67 diversity of the cluster. Consistent clusters and different clusters can be identified in each group.
68 This method can effectively identify the cluster of differences, and this theory is also used in our
69 method.

70 In order to improve the algorithm stability of MVCMO (Yu et al., 2019), we introduce the "Heat
71 Kernel" to measure similarity between patients. And we use Backtracking Line Search to find the
72 optimal solution more accurately. In this study, we propose a Multi-view Clustering based on Stiefel
73 Manifold (MCSM) method for multi-view clustering problems with potential clusters. Firstly, we
74 introduce a "Heat Kernel" to measure similarity. The patient-patient similarity network is
75 constructed using k-nearest neighbor (KNN) method. Then we establish a binary optimization
76 model for the simultaneous clustering problem. The solving process of the objective optimization
77 problem is divided into three steps. First, we project our target function onto each of Stiefel
78 manifold's tangent vector Spaces. Second, we do Backtracking Line Search on Stiefel manifold for
79 the objective problem. Third, we retract the found points to the manifold with singular value
80 decomposition. Finally, the KNN method is used for integrating the obtained clusters from three
81 omics to get the final result. The proposed MCSM method has two highlights. One is that it
82 preserves as much data information of each sample as possible. The other is that it can identify the
83 cluster effectively when the underlying clusters are different. We experiment on simulated datasets
84 to see the algorithm's performance when there are potential clusters. The experimental results on
85 simulated datasets and several multiple omics datasets from TCGA show that our method has better
86 performance than state-of-art methods.

87 **2.Datasets and Methods**

88 The overall design of our method is illustrated in Figure 1.



89

90 **Figure 1:** The process of MCSM method.

91 2.1 Datasets and preprocessing

92 In this paper, we selected four cancer datasets in the TCGA for experiment, including gene
93 expression data, miRNA expression data and DNA methylation data from samples of cancer patients.

94 The cancer datasets include glioblastoma multiforme (GBM) with 215 samples, breast invasive
95 carcinoma (BIC) with 105 samples, Skin Cutaneous Melanoma (SKCM) with 439 samples and
96 Acute Myeloid Leukemia (AML) with 96 samples.

97 Firstly, if the data of a patient loses more than 20% in any data type, the patient will be deleted.

98 Secondly, if the missing value of a feature in all patients exceeds 20%, it will be filtered out. Thirdly,
 99 the K-nearest-neighbor method is adopted to fill in missing data.

100 Fourthly, we log transform the data set to make it more stable. Finally, each feature is normalized
 101 in the constructed network to make it have a standard normal distribution. We performed the
 102 following normalization for each data type:

$$103 \quad \hat{f} = \frac{f - E(f)}{\sqrt{\text{Var}(f)}}. \quad (1)$$

104 where f is the characteristic of sample data, \hat{f} is the corresponding characteristic after f
 105 normalization of f , $E(f)$ and $\text{Var}(f)$ represent the sample mean and sample variance respectively.

106 2.2 Construction of the patient-to-patient similarity graph

107 Denoted $\{X^m\}_{m=1}^M$ as multi-view data from N patient samples, which has m data type in total.
 108 Each X^m is a matrix of $p_m \times N$, then a similar network graph G^m is constructed to reflect the
 109 neighborhood relationship between the samples.

110 In the similar network of type m , $G^m = (V^m, E^m, W^m)$, V^m is vertex set, E^m is edge set, and
 111 W^m is adjacency matrix. The adjacency matrix of W^m in graph G^m is a symmetric matrix.

112 In this paper, ‘‘Heat Kernel’’ is used to measure the similarity between samples (Ding et al., 2018).

$$113 \quad S_{ij}^m = \exp\left(-\frac{\|x_i^m - x_j^m\|^2}{2t^2}\right), i = 1 \dots, N, j = 1 \dots, N. \quad (2)$$

114 Next, we construct the K-nearest neighbor graph based on the similarity matrix S^m . If the vertex
 115 has an edge between v_i and v_j , then W_{ij}^m represents the edge weight, otherwise 0.

$$116 \quad W_{ij}^m = \begin{cases} S_{ij}^m, & v_j \in N_i, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

117 here, N_i is the neighborhood of v_i (including v_i), N_i with size k , and the number of k usually
 118 depends on the size of the sample. Essentially, we assume that local similarity is more reliable than
 119 remote similarity. This is a modest assumption, and it is widely used by other manifold learning
 120 algorithms (Ding et al., 2018).

121 2.3 Construction of objective optimize problem

122 The objective optimize problem of the spectral clustering method is:

$$123 \quad \min_{U_m \in \mathbb{R}^{N \times k}} \text{trace}(U_m^T L_m U_m)$$

$$124 \quad \text{s. t. } U_m^T U_m = I_K. \quad (4)$$

125 Here, the $L_m = (D_m - A_m)$. The A_m is the corresponding adjacency matrix of similar network
 126 G^m , and D_m is the diagonal matrix constructed using the degree of all the nodes in the m -th network.
 127 Then, used U_m for K-means clustering and find its minimum k eigenvectors in order to obtain the
 128 clustering labels.

129 Based on the spectral clustering, (Zhang et al., 2015) proposed a multi-view network clustering
 130 method. Its objective optimize problem is:

$$131 \quad \min \sum_{m=1}^M \sum_{k=1}^K \frac{(S_{i,k}^m)^T (D_m - A_m) (S_{i,k}^m)}{(S_{i,k}^m)^T (S_{i,k}^m)} - \beta \sum_{l \neq m} \sum_{k=1}^K \frac{(S_{i,k}^m)^T (S_{i,k}^l)}{\|S_{i,k}^m\|_2 \|S_{i,k}^l\|_2}$$

$$132 \quad \text{s. t. } S_{i,k}^m \in \{0,1\}, i = 1 \dots, N; m = 1 \dots, M; k = 1 \dots, k;$$

$$133 \quad \sum_{k=1}^K S_{i,k}^m = 1, \text{ for } m = 1 \dots, M. \quad (5)$$

134 The binary optimization problem cannot be solved in polynomial time. So, the objective function
 135 of multi-view spectral clustering can be constructed as follows:

136

$$\min_{U_m \in \mathbb{R}^{N \times k}} \text{trace}(U^T L U)$$

137

$$\text{s. t. } U^T U = I_K \quad (6)$$

138

$$\text{where, } L = \begin{pmatrix} L_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & L_m \end{pmatrix} - \beta \begin{pmatrix} 0 & \cdots & I_n \\ \vdots & \ddots & \vdots \\ I_n & \cdots & 0 \end{pmatrix}, U = \begin{pmatrix} U_1 \\ U_2 \\ \vdots \\ U_3 \end{pmatrix}.$$

139

β is used to balance the weight parameters between the network and within the network. If we

140

have abundant prior knowledge, we can set it according to prior information. Otherwise, when

141

building a network, we can try to establish a connection at the same level (e.g. similar connection

142

densities) and set it directly to 1. In our experiment, we set $\beta = 1$ directly.

143

However, the optimization problem (6) combines the information of all networks together and

144

will loss the information in each network. The proposed MVSM method still follows the original

145

objective function of multi-view spectral clustering and the construction of Laplace matrix (Yu et

146

al., 2019).

147

The objective optimization problems to be solved are as follows:

148

$$\min_{U_m \in \mathbb{R}^{N \times k}} \text{trace}(U^T L U)$$

149

$$\text{s. t. } U_m^T U_m = I \quad (7)$$

150

When we set $U_m = \frac{S^m}{\|S^m\|_2}$, the objective function $U_m^T U_m = I_K$ substitute as $\sum s_{i,j}^N = 1$. It

151

transforms the constraints for each network into one equation.

152

2.3 The solution of objective optimize problem

153

To solve the objective function (7), we project it onto the Stiefel manifold and solve it by linear

154

search. The process is roughly divided into three steps.

155 First, we project the target function $\text{trace}(U^T LU)$ onto each of Stiefel manifold's tangent vector
 156 Spaces.

157 The tangent vector space of M is

$$158 \quad TM_m = \{U_m B + (I - U_m U_m^T)C : B = -B^T, C \in \mathbb{R}^{N \times k}\}. \quad (8)$$

159 here, each Stiefel manifold $M_m = U_m \in \mathbb{R}^{N \times k}$: s. t. $U_m^T U_m = I_k$.

160 So, the negative gradient of the objective function can be expressed as:

$$161 \quad Z = -\nabla \text{trace}(U^T LU) = -LU = (Z_1^T, Z_2^T, \dots, Z_m^T)^T. \quad (9)$$

162 Then, we search the next point along the direction η_m on each tangent vector space of the
 163 manifold. Where,

$$164 \quad \eta_m = Z_m - \frac{1}{2} U_m^m ((U_m^m)^T U_m^m + Z_m^T). \quad (10)$$

165 Second, we do Backtracking Line Search on Stiefel manifold for problem (9-10).

166 The purpose of line search is to find the smallest point of the target function in the search direction.
 167 However, it is time-consuming to find the accurate minimum point. The search direction is already
 168 approximate, so we just to find the minimum point approximation at a lower cost. Backtracking
 169 Line Search (BLS) is such a Line Search algorithm. The idea of the BLS algorithm is to set an initial
 170 step size α_0 in the search direction. Then, if the step size is too large, we reduce the step size until
 171 it is appropriate.

172 Backtracking Line Search in the negative gradient direction of the objective function is as follow:

$$173 \quad f(U + \alpha \eta) \leq f(U) + \alpha cm$$

$$174 \quad m = \eta^T Z \quad (11)$$

175 Where, η is the current search direction, α is the step size, and c is the control parameter,
 176 which needs to be manually verified according to the situation.

177 If the current U does not satisfy inequation (11), then a parameter τ is required to adjust the step
 178 size:

$$179 \quad \alpha = \tau\alpha \quad (12)$$

180 Where, the parameter τ is controls the reduce search step size.

181 Third, we retract the found points to the manifold with singular value decomposition.

$$182 \quad U = W\Sigma U^T, U = WU^T. \quad (13)$$

183 After the manifold optimization process, we get the values of U . The whole process of our
 184 proposed method is summarized in Algorithm 1.

185

Algorithm 1: The linear search based on Stiefel manifold

Step 1. The negative gradient direction of the objective function is projected onto Stiefel manifold.

$$\eta_m = Z_m - \frac{1}{2}U^m((U^m)^T U^m + Z_m^T);$$

$$\eta = (\eta_1^T, \eta_2^T, \dots, \eta_m^T)^T;$$

Step 2. Backtracking Linear search in tangent vector space:

$$U = U + \alpha\eta, \alpha \in (0,1);$$

$$f(U + \alpha\eta) \leq f(U) + \alpha cm;$$

$$\mathbf{m} = \eta_m^T \mathbf{Z};$$

Step 3. Retracted the points obtaining in 2 to the Stiefel manifold:

$$U = W\Sigma U^T, U = WU^T;$$

Step 4. Repeat 1- 3 until the convergence condition is satisfied:

$$((r_err_f > 1e - 8) \parallel (r_err_g > 1e - 4)) \&\& (iters < 1000);$$

Where r_err_f represents the relative error of objective function f , r_err_g represents the relative error of the negative gradient of the objective function Z , and $iters$ represents the number of iterations.

186 Here, we get the solution of the objective function, and then we perform k-means to cluster U and
187 obtain the cluster labels C_1, C_2, \dots, C_k . Finally, we integrate the clustering results obtained from
188 three omics by using k-nearest neighbor method.

189 Remark: we set $c = 6, \tau = 0.1$ in our experiment. We will set out the reasons in **3.1.2**.

190 **3.Results**

191 In this section, to show the effective of our proposed method, we compare it with AASC algorithm
192 (Huang et al., 2012), SNF method (Wang et al., 2014), MOCMO (Yu et al., 2019), MVSC (Zhang
193 et al., 2015) and Grassmann manifold clustering method (Ding et al., 2018).

194 **3.1The selection of parameter**

195 **3.1.1 The number of clustering**

196 When the clusters k is not known, we can select it according the value of silhouette. Firstly, we
197 did experiments with k equals 2 to 10. Then, we choose the number corresponding to the maximum
198 silhouette coefficient. To compare MVSM to other methods, we set k as a known value.

199 **3.1.2 The Backtracking Line Search parameters**

200 There are three parameters in the Backtracking Line Search parameters, η , α and c . Where, η
201 is the current search direction, α is the step size, and c is the control parameter, which needs to be
202 manually verified according to the situation. Firstly, we initialize $\alpha = 0.01$. During the experiment,
203 it was found that if the value of c was too small, the step size would not be adjusted during the
204 search process. However, if we want to adjust appropriately, then we need to set the parameter c
205 according to the objective function value and gradient value of the initial point. Therefore, according
206 to several data sets used in the experiment, we set $c = 6, \tau = 0.1$.

207 **3.2 Experimental results on simulated datasets**

208 We refer to the same method in the article (Yu et al., 2019) to verify that this method is suitable
209 for datasets with uneven distribution of underlying clusters. Here we use a simulated dataset.

210 Since these methods (SNF, AASC and MVSC) were proposed using network tools, we simulate
211 the network structure firstly. For M omics networks, given the number of nodes N , these nodes are
212 assigned to K clusters. Then, the connections within the same cluster and different clusters are
213 generated. The probability of connections within a given cluster is greater than the probability of
214 connections between clusters.

215 We used the following four connection probability matrices to define the connections within and
216 between clusters to see the results when the connections between clusters changed:

217
$$P_1 = \frac{1}{N} \begin{pmatrix} 16 & 0 & 0 \\ 0 & 18 & 0 \\ 0 & 0 & 17 \end{pmatrix}, \quad P_2 = \frac{1}{N} \begin{pmatrix} 16 & 0.4 & 0.6 \\ 0.4 & 18 & 0.55 \\ 0.6 & 0.55 & 17 \end{pmatrix},$$

218
$$P_3 = \frac{1}{N} \begin{pmatrix} 16 & 0.8 & 1.2 \\ 0.8 & 18 & 1.1 \\ 1.2 & 1.1 & 17 \end{pmatrix}, \quad P_4 = \frac{1}{N} \begin{pmatrix} 16 & 1.2 & 1.8 \\ 1.2 & 18 & 1.65 \\ 1.8 & 1.65 & 17 \end{pmatrix}.$$

219 Where, each term (i, i) of the four matrices defines the connection probability within the cluster,
 220 and each term $(i, j), i \neq j$ represents the connection probability between cluster i and cluster j . For
 221 each setting, we consider that in all M omics, the sample size of each category is the same, and the
 222 omics of distribution is different. To see the performance of the method, we tested two settings:

223 Setting 1: $N=150, M=3$, cluster size:(50,50,50); (30,90,30); (40,60,50);

224 Setting 2: $N=1000, M=6$, cluster size:(300,300,400); (300,300,400); (400,300,300); (300,350,350);
 225 (300,400,300); (450,250,300).

226 We use the Rand index to evaluate the clustering performance, which is defined as:

227
$$RI = \frac{TP + YN}{TP + FP + TN + FN},$$

228 with the 'TP', 'TN', 'FP', 'FN' are present the number of the true positive, true negative, false positive,
 229 and false negative. 'TP' is defined as the number of intersection nodes in the same cluster, which are
 230 also clustered in the same cluster, and other nodes are defined similarly.

231 On this basis, we obtain the rand index comparison of several methods:

232 **Table 1:** Performance comparison of different methods. $N=150, M=3$,
 233 cluster size:(50,50,50); (30,90,30); (40,60,50)

Method	P_1	P_2	P_3	P_4
--------	-------	-------	-------	-------

AASC	0.73	0.73	0.73	0.73
SNF	0.68	0.67	0.67	0.67
MVSC	0.99	0.98	0.94	0.87
MCSM	1	0.99	0.94	0.97

234

235 **Table 2:** Performance comparison of different methods. N=1000, M=6, cluster size:(300,300,400),
 236 (300,300,400); (400,300,300); (300,350,350); (300,400,300); (450,250,300)

Method	P₁	P₂	P₃	P₄
AASC	0.75	0.75	0.75	0.75
SNF	0.75	0.75	0.75	0.75
MVSC	0.93	0.93	0.93	0.93
MCSM	0.96	0.95	0.95	0.95

237

238 For each setting, we run it 50 times and take the average of the results. From Table 1 and Table
 239 2, we can see that all four methods cluster nodes with an average Rand index is close to 1 when the
 240 cluster sizes of different groups are the same. This also illustrates the importance of data integration.
 241 When the size of the underlying cluster is different, because both AASC and SNF set the underlying
 242 cluster to be the same, they cannot detect the difference between the different views, so the MVSC
 243 and our method have better performance. Table 1 and Table 2 show the mean Rand index and its

244 standard deviation in the bracket when the underlying clusters are different for the two setups. By
 245 using more strict relaxation of the binary variables, more information of the clusters can be kept.
 246 our method out performs MVSC in both setups and when there are less networks, our method
 247 performs more stable, which achieves much less standard deviation of the Rand index.

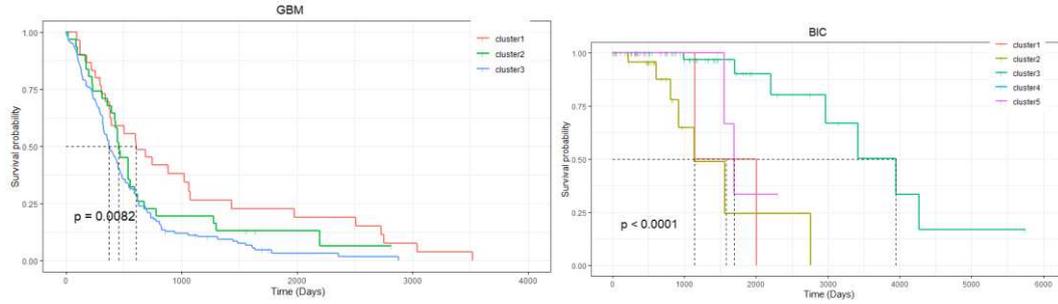
248 3.3 Experimental results on real datasets

249 In order to prove the effectiveness of our method on the real datasets. We apply our method on
 250 multiple datasets analyzed by (Wang et al., 2014), and compare our results with those of SNF fused
 251 by similarity network. The final Cox survival P values are listed in Table 3.

252 **Table 3:** Comparison of Cox survival p-values.

Cancer type	SNF	Grassmann Cluster	MOCMO	AASC	MVSC	MCSM
GBM(3)	0.0002	0.0043	0.0019	0.0022	0.00072	0.0001
BIC(5)	0.0011	0.0002	0.00016	0.00015	0.0007	0.0025
SKCM(4)	0.0001	0.19	0.00045	0.00016	0.00045	0.0001
AML(5)	0.037	0.12	0.03	0.045	0.058	0.019

253 Integrative clustering discovers clinically significant subtypes of cancer. Shown are Kaplan Meier
 254 plots of the overall survival of integrative clusters for Glioblastoma Multiforme (GBM) (a), breast
 255 Invasive carcinoma (BIC) (b), Skim Cutaneous Melanoma (SKCM) (c) and Acute Myeloid
 256 Leukemia (AML) (d) in figure 2. P-values are computed from the log rank test.

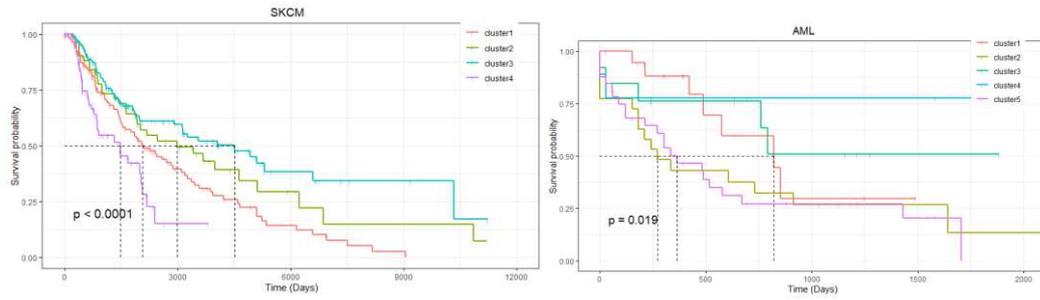


257

258

(a)

(b)



259

260

(c)

(d)

261

Figure 2: Survival plots for GBM, BIC, SKCM, and AML tumors.

262

It can be seen from the Table 3, in three of the four cancers datasets (GBM, BIC, SKCM and AML), our method obtains more significant differences than comparison methods in survival time.

263

264

For BIC dataset, the insignificant difference may be due to the small cluster difference of the data

265

itself. Survival plots for GBM, BIC, SKCM, and AML tumors are shown in Figure 2. Our proposed

266

method can also be extended to improve the survival rate prediction task.

267

3.4 A case study: Comparison of clusterings to established subtypes

268

In order to compare the results of our clustering with the established biological subtypes, we

269

downloaded the clinical data of 215 GBMs from the cBio Cancer Genomis Portal

270

(<http://www.cbioportal.org/>) at the Memorial Sloan-Kettering Cancer Center. For the GBM, there

271 exist four established subtypes determined by patients' gene expression profiles, which are Classical,
 272 Mesenchymal, Neural and Proneural (Verhaak et al., 2010), and a subtype called Glioma-CpG island
 273 methylator phenotype (G-CIMP) generated by DNA methylation clustering (Noushmehr et al.,
 274 2010). Comparison of our GBM clustering with these existing subtypes (Table 4) shows that our
 275 method not only reflects evidence from one data type, but also finds a cluster that considers both
 276 gene expression and DNA methylation information.

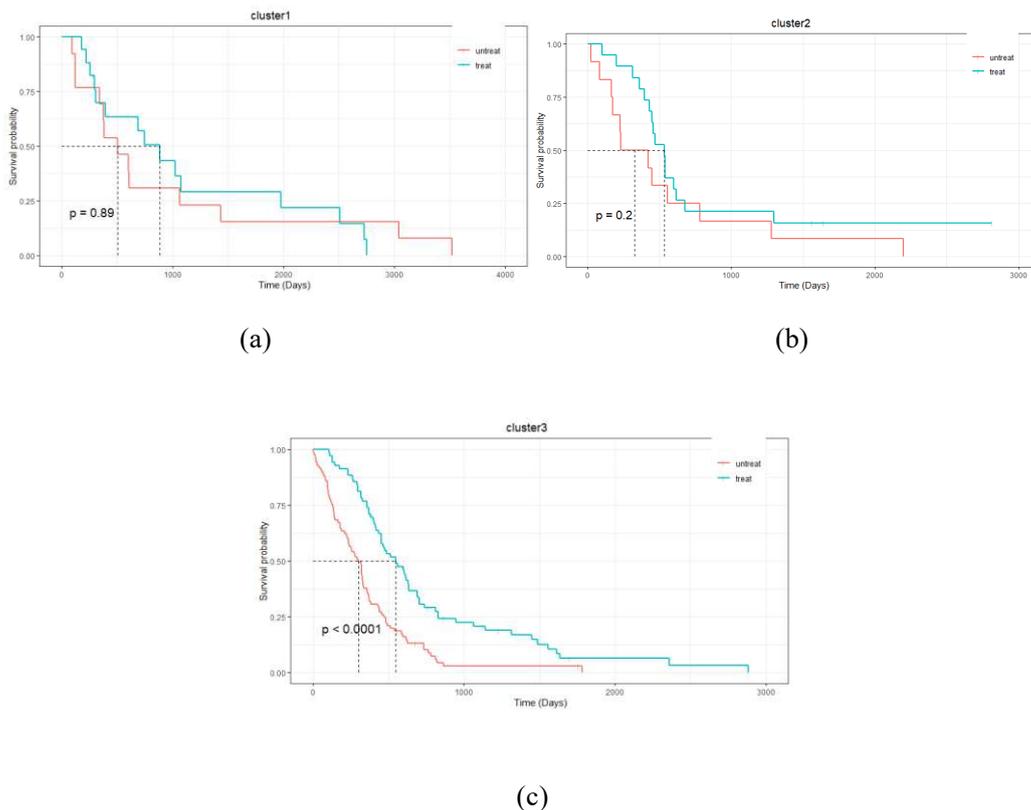
277 **Table 4:** Comparison of clusterings to established subtypes

Clusters	Gene expression subtypes (Verhaak et al., 2010)				DNA methylation subtypes (Noushmehr et al., 2010)	
	Classical	Mesenchymal	Neural	Proneural	G-CIMP	Non-G-CIMP
Cluster1	0	0	1	13	16	14
Cluster2	2	25	2	2	0	31
Cluster3	55	41	31	24	2	152

278

279 For gene expression subtypes, it can be seen that cluster 1 mainly contains Proneural subtype,
 280 cluster 2 mainly contains Proneural subtype, and they have strong enrichment. However, for DNA
 281 methylation subtypes, G-CIMP subtypes are mainly distributed in cluster 1. If only the DNA
 282 methylation information is considered, cluster2 and cluster3 are likely to merge. So, we can
 283 conclude that it's important to consider both gene expression and DNA methylation information.

284 In order to further understand the biological significance of clusters, we investigated the response
285 to temozolomide (TMZ) treatment of the GBMs. TMZ is an alkylation agent that causes incorrect
286 pairing of thymine during DNA replication. In the GBM dataset, 105 patients were treated with
287 TMZ. Figure 3 indicated that the TMZ-treated samples had different drug responses compared to
288 the samples not treated with the drug. For different clusters, the degree of drug response of TMZ
289 was also different. Compared with Cluster 1 and Cluster 2, patients in Cluster 3 had significantly
290 increased survival time after treatment with TMZ (P value using Cox log-rank test=0.0001), and
291 this medication was also more meaningful. The results show that the clusters we obtained can be
292 used as a reference for identifying the effectiveness of drugs.



297 **Figure 3:** Survival analysis of GBM patients for treatment with Temozolomide in the different
298 clusterings.

299 **4. Conclusion**

300 Multi-view data clustering is a hot topic in recent years. Recent work has focused on cases where
301 the underlying clusters are consistent, and as we reviewed in the first section, several approaches
302 have been proposed. When the underlying cluster is different, some methods are proposed to find
303 different clusters. However, as we know, both consistent and differentiated clusters can exist at the
304 same time. This leads us to study multi-view simultaneous clustering to find both consistent and
305 different cluster data. In this paper, we propose a multi-view clustering model. On the basis of
306 manifold optimization, the algorithm for formula optimization is proposed. Simulation results show
307 that the performance of the proposed method is better than that of the existing algorithm under the
308 same underlying cluster condition. We download the gene expression, miRNA expression and DNA
309 methylation datasets of GBM, BIC, SKCM and AML from TCGA, and also carry out numerical
310 experiments, showing that our method is superior to several comparison methods. In the future work,
311 the cluster difference problem is still worth researching, and we will integrate other omics
312 information such as gene mutation data.

313 **Declarations**

314 **Ethics approval and consent to participate:** Not applicable.

315 **Consent for publication:** Not applicable.

316 **Availability of data and materials:**

317 All data generated or analysed during this study are included in this published article [and its suppl
318 ementary information files]

319 **Competing interests:** Not applicable.

320 **Funding:** Not applicable.

321 **Authors' contributions:** T J' contributions are processed a multi-view clustering based on Stiefel

322 manifold method (MCSM) and the numerical simulation results of MCSM model. Z JP'
323 contribution lies in the embellishment of the article. Z CH' contribution lies in the thought
324 guidance of the method.

325 **Acknowledgements:** This work was supported by grants from the Xinjiang Autonomous Region
326 University Research Program (No. XJEDU2019Y002), and the National Natural Science
327 Foundation of China (No. U19A2064, 61873001).

328 **References**

- 329 [1] Absil PA, Mahony R and Sepulchre R. Optimization Algorithms on Matrix Manifolds.
330 Princeton University Press, Princeton, NJ (2008).11462-11467. doi: 10.1515/9781400830244
- 331 [2] Alon U, Barkai N, Notterman DA, Gish K, Ybarra S, Mack D, and Levine AJ. Broad patterns
332 of gene expression revealed by clustering analysis of tumor and normal colon tissues probed
333 by oligonucleotide arrays. Proceedings of the National Academy of Sciences of the United
334 States of America (1999) 96(12): 6745-6750. doi: 10.1073 / pnas.96.12.6745
- 335 [3] Ashburner M, Ball CA, Botstein D and Blake JA. Gene ontology: tool for the unification of
336 biology. Nature Genetics (2000) 25(1): 25-29. doi: 10.1038/75556
- 337 [4] Bendor A, Shamir R and Yakhini Z. Clustering gene expression patterns. Journal of
338 Computational Biology (1999) 6: 281-297. doi: 10.1145/299432.299448
- 339 [5] Bickel PJ and Chen A. A nonparametric view of network models and new managirvan and
340 other modularities. Proceedings of the National Academy of Sciences of the United States of
341 America (2009)106(50): 21068-21073. doi: 10.1073/pnas.0907096106

- 342 [6] Cancer Genome Atlas Research Network. Comprehensive genomic characterization of
343 squamous cell lung cancers. *Nature* (2012) 489: 519-525. doi: 10.1016/j.ypat.2012.11.048
- 344 [7] Chang HY, Nuyten DSA, Sneddon JB, Hastie Tibshirani R, Dai HY and He YD. Robustness,
345 scalability and integration of a wound-response gene expression signature in predicting breast
346 cancer survival. *Proceedings of the National Academy of Sciences of the United States of*
347 *America* (2005) 102(10): 3738-3743. doi:10.1073/pnas.0409462102
- 348 [8] Chalise P, Koestler DC, Bimall M, Yu Q. Integrative clustering methods for high dimensional
349 molecular data. *Translational cancer research* (2014) 3(3), 202–216. doi: 10.3978/j.issn.2218-
350 676X.2014.06.03
- 351 [9] Chen C, Ng MK, and Zhang S. Block spectral clustering methods for multiple graphs.
352 *Numerical Linear Algebra with Applications* (2017) 24(1): 1-20. doi: 10.1002/nla.2075
- 353 [10] Dai LY, Feng CM, Liu JX, Zheng CH, Yu J. Robust graph regularized discriminative
354 nonnegative matrix factorization for characteristic gene selection. *IEEE International*
355 *Conference on Bioinformatics and Biomedicine* (2016) 1253-1258. doi:
356 10.1109/BIBM.2016.7822698
- 357 [11] Ding H, Michael S, Wang C. Integrative cancer patient stratification via subspace merging.
358 *Bioinformatics* (2018) 35(10): 1653-1659. doi: 10.1093/bioinformatics/bty866
- 359 [12] Fridley BL, Lund SP, Jenkins GD, Wang LW. A Bayesian integrative genomic model for
360 pathway analysis of complex traits. *Genetic epidemiology* (2012) 36:352-359. doi:
361 10.1002/gepi.21628

- 362 [13] Ge SG, Xia J, Sha W, Zheng CH. Cancer subtype discovery based on integrative model of
363 multigenomic data. *IEEE/ACM transactions on computational biology and bioinformatics*
364 (2017) 14(5): 1115-1121. doi: 10.1109/TCBB.2016.2621769
- 365 [14] Ha MJ, Baladandayuthapani V and Do KA. Dingo: Differential network analysis in genomics.
366 *Bioinformatics* (2015) 31(21): 3413-3420. doi: 10.1093/bioinformatics/btv406
- 367 [15] Noushmehr H, Weisenberger DJ, Diefes K, Phillips HS, Pujara K, Berman BP and Cancer
368 Genome Atlas Research Network. Identification of a CpG island methylator phenotype that
369 defines a distinct subgroup of glioma. *Cancer cell* (2010) 17(5): 510-522.
- 370 [16] Huang HC, Chuang YY and Chen CS. Affinity aggregation for spectral clustering. *Conference*
371 *on Computer Vision and Pattern Recognition* (2012) 2012: 773-780. doi:
372 10.1109/CVPR.2012.6247748
- 373 [17] Huang BY, Zhang HB, Gu LJ, Ye BX and Xiong X. Advances in immunotherapy for
374 glioblastoma multiforme. *Clinical Developmental Immunology* (2017) 2017(3): 1-11. doi:
375 10.1155/2017/3597613
- 376 [18] Huang E, Cheng SH, Dressman H and Pittman J. Gene expression predictors of breast cancer
377 outcomes. *The Lancet* (2003) 361(9369): 1590-1596. doi: 10.1016/S0140-6736(03)13308-9
- 378 [19] Kanehisa M, Goto S, Furumichi M, Tanabe M and Hiraakawa M. Kegg for representation and
379 analysis of molecular networks involving diseases and drugs. *Nucleic Acids Research* (2010)
380 38(suppl_1): D355-D360. doi: 10.1142/9781848165632_0020
- 381 [20] Kumar R, Kamdar D, Madden L, Hills C. Th1/th2 cytokine imbalance in meningioma,

382 anaplastic astrocytoma and glioblastoma multiforme patients. *Oncology Reports* (2006) 15(6):
383 1513-1516. doi: 10.3892/or.15.6.1513

384 [21] Shen R, Olshen AB and Ladanyi M. Integrative clustering of multiple genomic data types using
385 a joint latent variable model with application to breast and lung cancer subtype analysis.
386 *Bioinformatics* (2010) 25(22): 2906-2912. doi: 10.1093/bioinformatics/btp659

387 [22] Shen R, Mo Q, Schultz N and Seshan VE. Integrative subtype discovery in glioblastoma using
388 icluster. *Plosone* (2012) 7(4): e35236. doi: 10.1371/journal.pone.0035236

389 [23] Speicher NK and Pfeifer N. Integrating different data types by regularized unsupervised
390 multiple kernel learning with application to cancer subtype discovery *Bioinformatics* (2015)
391 31(12): i268-i275. doi: 10.1093/bioinformatics/btv244

392 [24] Stefansson OA, Moran S, Gomez A and Sayols S. A DNA methylation-based definition of
393 biologically distinct breast cancer subtypes. *Molecular Oncology* (2015) 9(3): 555-568. doi:
394 10.1016/j.molonc.2014.10.012

395 [25] Verhaak RGW, Hoadley KA, Purdom E, Wang V, Qi Y, and Wilkerson MD. Integrated genomic
396 analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities
397 in *pdgfra*, *idh1*, *egfr*, and *nfl*. *Cancer Cell* (2010) 17(1): 98-110.

398 [26] West M, Blanchette C, Dressman H and Huang E. Predicting the clinical status of human breast
399 cancer by using gene expression profiles. *Proceedings of the National Academy of Sciences of*
400 *the United States of America* (2001) 98(20): 11462-11467. doi: 10.1073/pnas.201162998

401 [27] Wang B, Mezlini AM, Demir F and Fiume M. Similarity network fusion for aggregating data

402 types on a genomic scale. *Nature Methods* (2014) 11(3): 333-337. doi: 10.1038/nmeth.2810

403 [28] Yuan Y, Savage RS, Markowetz F. Patient-specific data fusion defines prognostic cancer
404 subtypes. *Plos Computational Biology* (2011) 7(10): e1002227. doi:
405 10.1371/journal.pcbi.1002227

406 [29] Yu Y, Zhang LH and Zhang SQ. Simultaneous clusterin of multiview biomedical data using
407 manifold optimization, *Bioinformatics* (2019) 35(20): 4029-4037. doi:
408 10.1093/bioinformatics/btz217

409 [30] Zhang D, Chen P, Zheng CH, and Xia JF. Identification of ovarian cancer subtype-specific
410 network modules and candidate drivers through an integrative genomics approach. *Oncotarget*
411 (2016) 7(4): 4298. doi: 10.18632/oncotarget.6774

412 [31] Zhang S, Zhao H and Ng MK. Functional module analysis for gene coexpression networks
413 with network integration. *IEEE/ACM Transactions on Computational Biology and*
414 *Bioinformatics* (2015) 12(5): 1146-1160. doi: 10.1109/TCBB.2015.2396073

415 [32] Zheng CH, Ng TY, Zhang L, Shiu, C. K., & Wang, H. Q. Tumor classification based on non-
416 negative matrix factorization using gene expression data. *IEEE transactions on nanobioscience*
417 (2011) 10(2): 86-93. doi: 10.1109/TNB.2011.2144998

418 [33] Zheng CH, Yang W, Chong YW, Xia JF. Identification of mutated driver pathways in cancer
419 using a multi-objective optimization model. *Computers in Biology and Medicine* (2016) 72:
420 22-29. doi: 10.1016/j.combiomed.2016.03.002

Figures

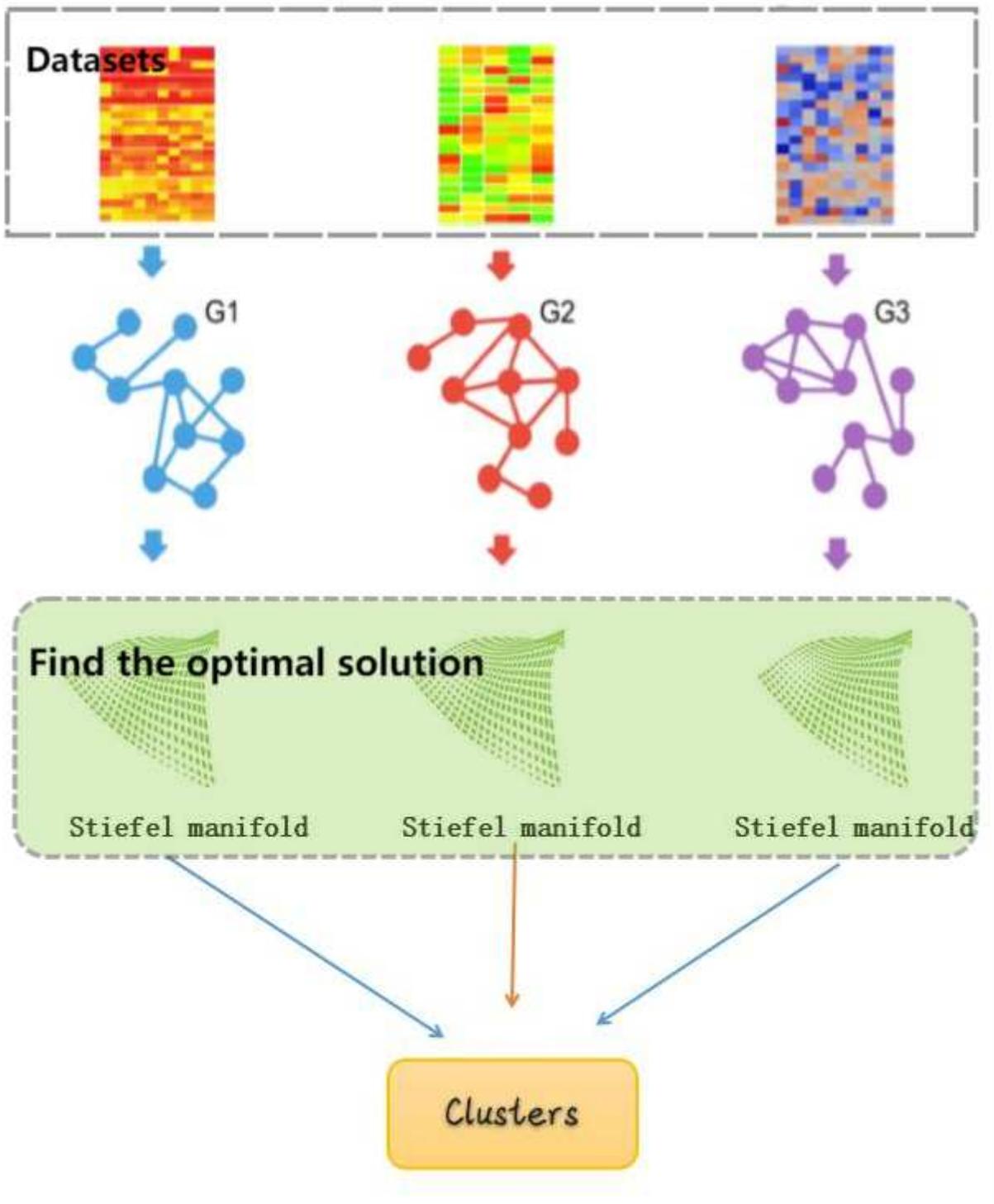
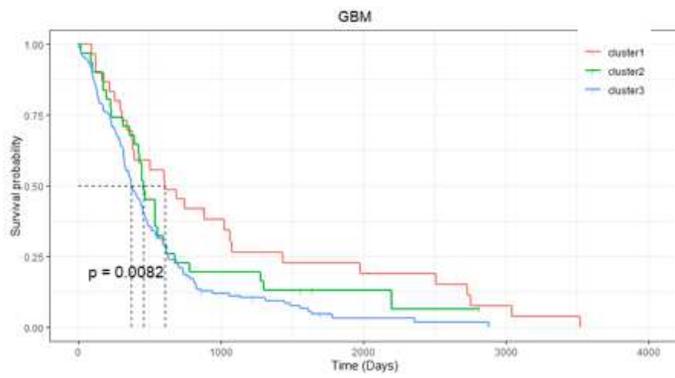
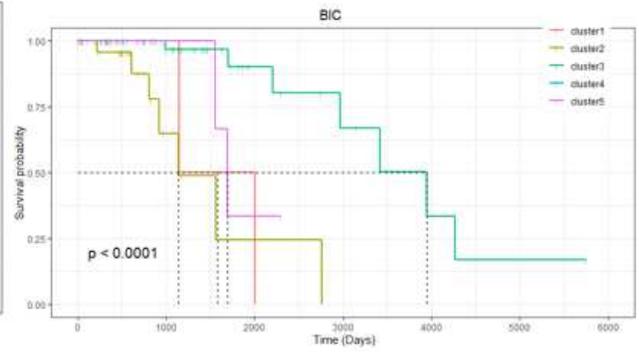


Figure 1

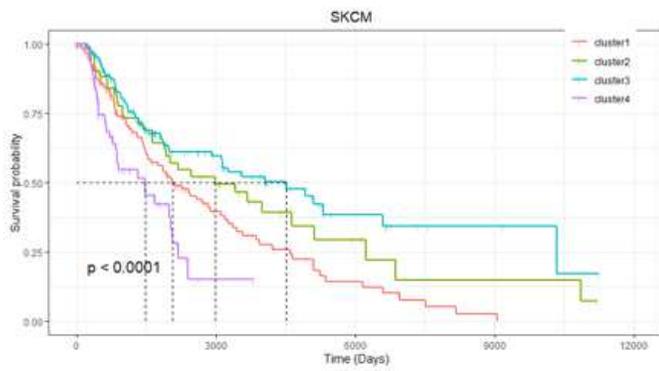
The process of MCSM method.



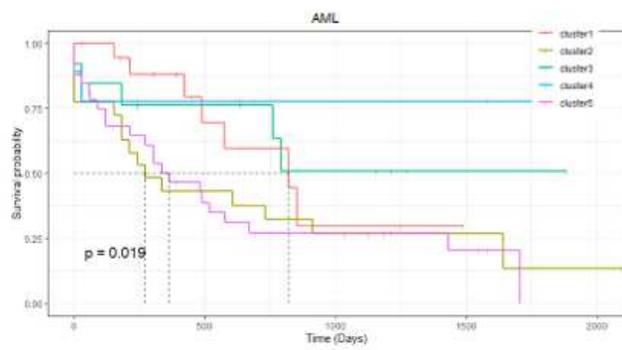
(a)



(b)



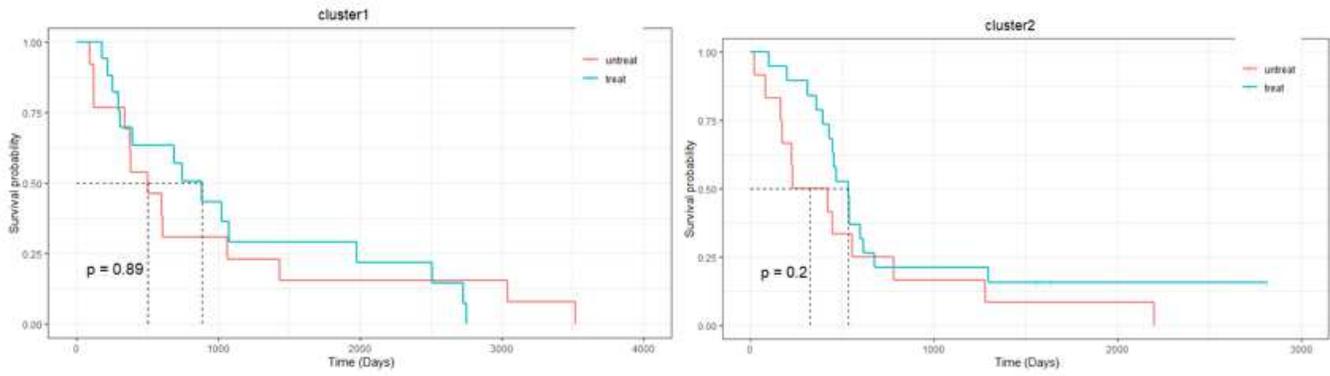
(c)



(d)

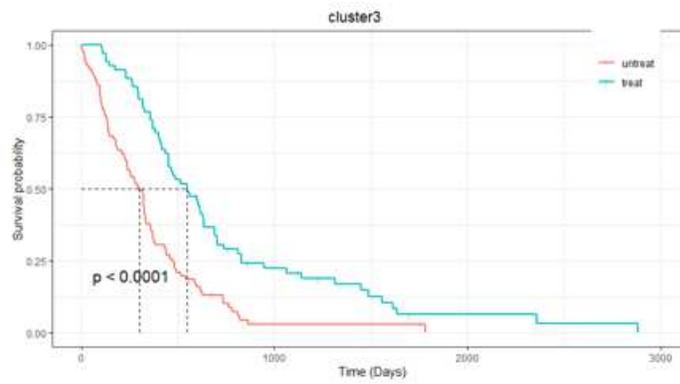
Figure 2

Survival plots for GBM, BIC, SKCM, and AML tumors.



(a)

(b)



(c)

Figure 3

Survival analysis of GBM patients for treatment with Temozolomide in the different clusterings.