

Performance Evaluation of Consensus Clustering Algorithms

STAT 598 Progress Report

Derek Chiu

August 25, 2015

Contents

1	Preface	2
2	Introduction	2
3	Methods	2
3.1	Clustering Algorithms	2
3.1.1	Hierarchical Clustering	2
3.1.2	k-Means Clustering	3
3.1.3	k-Medoids Clustering	3
3.1.4	Nonnegative Matrix Factorization	3
3.2	Consensus Clustering	3
4	External Evaluation	5
4.1	Purity and Entropy	5
4.2	Kappa Statistics	5
4.3	Adjusted Rand Index	5
5	Internal Evaluation	6
5.1	Proportion of Ambiguously Clustered Pairs (PAC)	6
5.2	Davies-Bouldin Index	7
5.3	Dunn Index	7
5.4	Rousseeuw's Silhouette	7
5.5	Silhouette Average Width	8
5.6	C-Index	8
5.7	Gamma Index	8
5.8	CH Index	9
5.9	Summary	9
6	Ranked Indices	9

7	Simulations	10
7.1	Worms Dataset	10
7.2	Rings Dataset	13
8	References	14

1 Preface

This progress report fulfills the UBC Science Co-op requirement to submit a work term report at the end of every four month period. BC Cancer Agency (BCCA) is a not-for-profit organization that aims to provide care for cancer patients and conduct innovative cancer research. Our department, OvCaRe, is the Ovarian Cancer Research team tasked with studying ovarian cancers of many types. The objective of the project I am working on is to discover a viable classifier for ovarian high-grade serous carcinoma (HGSC). My role is to help devise a clustering algorithm that can statistically partition subtypes of HGSC without knowledge of the pathological properties of each sample. The progress report will evaluate the method we are using, consensus clustering, on a publicly available data set as well as simulated data sets. The final technical report will contain results of our method applied on HGSC data from our own cohort.

2 Introduction

Unsupervised learning is the process of inferring something about a data structure without knowing its true class labels. Cluster analysis is an unsupervised learning method of assigning entities into different groups based on one or more of their attributes. It is unsupervised because we do not know the true partitions of the entities. The objective is to place similar objects together in the same cluster and separate dissimilar objects into different clusters. For example, in genomics studies, we frequently try and cluster patient samples measured on a large number of molecular features. When we get a clustering assignment from an algorithm, we often want to evaluate its performance. Ideally, a good clustering algorithm is able to differentiate entities without knowledge of the true class labels. In addition, we want the algorithm to arrive at a stable and optimal number of clusters. The choice of the number of clusters is not trivial in some cases.

3 Methods

3.1 Clustering Algorithms

There are many clustering algorithms, each approaching the clustering problem in a different way. It is most important to note the advantages and limitations of each algorithm. There are some definitions to take note of:

- **Compactness:** how close the objects are in each cluster
- **Connectivity:** how connected the objects are in the feature space

3.1.1 Hierarchical Clustering

This clustering algorithm is very popular because of its intuitive representation using dendrograms (trees). Based on a distance matrix, the objects/features are clustered based on a linkage type. More similar objects are joined near the bottom of a dendrogram whereas less similar objects are joined at a higher tree height. A linkage criterion determines the distances amongst a set of objects/features using the pairwise distances. For

example, an average linkage would use the average pairwise distances. In this way, a dendrogram with all objects/features can be made by recursively linking larger and larger sets of observations together.

CITATION: it turns out that single linkage works very well for data sets exhibiting connectivity but not compactness. An example of this would be a tree rings. Clusters are circles, and objects that are far away are in the same cluster compared to objects that are close. On the other hand, average linkage works well for data sets exhibiting compactness.

3.1.2 k-Means Clustering

First, k means are randomly initialized in the multi-dimensional object/feature space that we wish to cluster. Assigning each object/feature to its closest mean forms the clusters. The k means are re-calculated based on the centroids (center points) within each cluster. This process is repeated until the centroids converge.

There are two caveats to note when using k -means. First, the cluster assignments are unstable because they depend on the random initialization of the means. We would preferably want to repeat this process many times to see whether the clusters change drastically. Secondly, the choice of k is not arbitrary. Cross-validation using an appropriate loss function is a popular method for choosing k .

3.1.3 k-Medoids Clustering

Very similar to k -means except that we initialize a set of random data points as medoids instead of random means that are generally not real data points. The most common version of k -medoids is PAM.

3.1.4 Nonnegative Matrix Factorization

Given a non-negative data matrix A , we can factor it into two matrices W and H , which are also non negative. W and H have important properties. Suppose A has genes as rows and samples as columns. If we are interested in clustering samples, then H has a reduced gene space of meta genes that fully explain the samples. Samples are clustered based on the metagene they are most associated with. If we are interested in clustering genes, then W has a reduced sample space of metasamples that fully explain the genes. Genes are clustered based on the metasample they are most associated with.

In gene expression data, it is common to standardize the genes. Doing so would disrupt the nonnegativity of A required for NMF. A simple remedy can solve this problem. We append the matrix $-A$ to the bottom of A , preserving the same number of columns, and set all the negative entries to 0. The computational complexity has been doubled as a result. NMF takes a long time to run, but offers the most stable clustering assignments.

3.2 Consensus Clustering

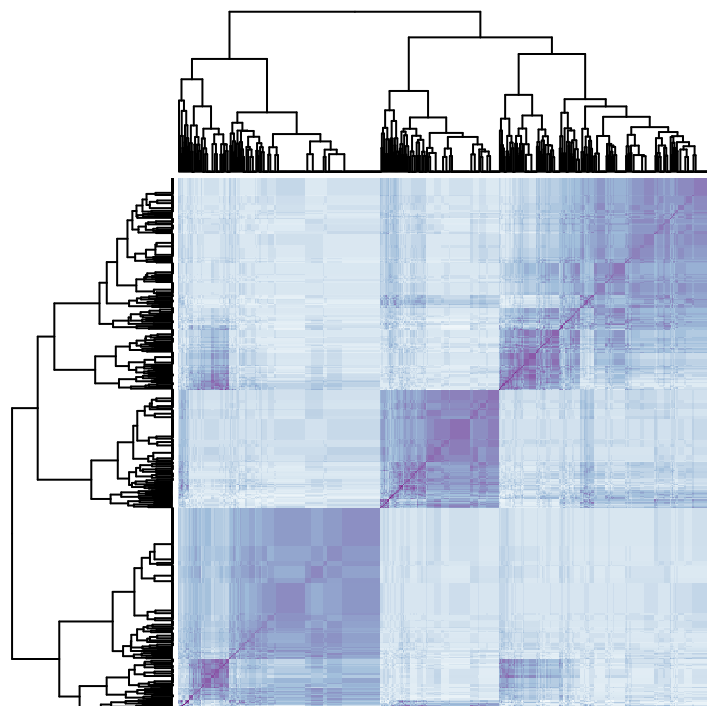
Consensus clustering is as much an algorithm itself as it is a way of combining realizations of other clustering algorithms. For this, it deserves its own section. Algorithms like k -means and PAM are unstable, as the clustering assignments depend on the initialization of centroids and medoids, respectively. Consensus clustering compares results from repeated runs of a clustering algorithm. Typically, these runs use different subsamples of the data to model sampling variability. A final clustering assignment is determined based on agreement across replicates.

A consensus matrix is a significant aspect of consensus clustering. All entries range from 0 to 1, and represent the proportion of replicates in which two objects/features were clustered together. A perfect clustering would consist of a consensus matrix with only 0s and 1s. More importantly, we can perform hierarchical clustering on a consensus matrix, resulting in a heat map with a diagonal block structure.

Repeating this process for different values of k (number of clusters) and looking at the corresponding heatmaps, we can see which value of k provides the most stable clustering assignment. Certainly, using a performance

measure such as PAC would be a more formal method of assessment that doesn't rely on potentially subjective visualizations.

Consensus clustering attempts to stabilize results from randomness introduced by subsampling and the clustering algorithm. A natural extension is sometimes called meta consensus clustering, where we want to compare results across different algorithms. For instance, we may combine the consensus clustering assignments from 10 different algorithms and come up with a final clustering.



The proportion of cases with at least 0.6 agreement is 0.2170958.

The confusion matrix is shown below, as well as different metrics for each class.

	C1	C2	C4	C5
C1	104	10	18	7
C2	1	37	10	1
C4	1	60	91	21
C5	3	0	16	109

	Sensitivity	Specificity	Pos Pred Value	Neg Pred Value	Prevalence	Detection Rate	Detection Prevalence	Balanced Accuracy
Class: C1	0.9541	0.9079	0.7482	0.9857	0.2229	0.2127	0.2843	0.931
Class: C2	0.3458	0.9686	0.7551	0.8409	0.2188	0.07566	0.1002	0.6572
Class: C4	0.6741	0.7684	0.526	0.8608	0.2761	0.1861	0.3538	0.7212
Class: C5	0.7899	0.9459	0.8516	0.9197	0.2822	0.2229	0.2618	0.8679

4 External Evaluation

External evaluation usually refers to the case when we compare our clustering assignments to true class labels, or have some gold standard to compare to. In applications, this might be the published clustering result. The downside of using external evaluation is that the reference classes may not be correctly clustered themselves, and we are treating these as the norm. None the less, we can explore a few metrics.

4.1 Purity and Entropy

Purity is defined as the sum of the entities of maximal class in each cluster divided by the total number of entities. The equation is:

$$Purity = \frac{1}{n} \sum_{r=1}^k \max_i(n_r^i)$$

where n is the total number of entities, k is the number of clusters, i is a particular class, and n_r^i is the number of objects classified into class i in cluster r . The larger the purity, the better the clustering accuracy.

##	NMF (Divergence)	NMF (Euclidean)	PAM (Spearman)	KM (Spearman)
##	0.7668712	0.7443763	0.7361963	0.7157464
##	HC (Diana)	PAM (Euclidean)	HC (Euclidean)	KM (Euclidean)
##	0.7116564	0.6687117	0.6237219	0.5971370
##	KM (MI)	PAM (MI)		
##	0.4376278	0.4130879		

Entropy measures the amount of uncertainty in each cluster. The equation is:

$$Entropy = -\frac{1}{n \log q} \sum_{r=1}^k \sum_{i=1}^q n_r^i \log \frac{n_r^i}{n_r}$$

The smaller the entropy, the less uncertain we are of the cluster membership, and the better the clustering performance.

##	NMF (Divergence)	NMF (Euclidean)	PAM (Spearman)	HC (Diana)
##	0.5101189	0.5322451	0.5470938	0.5639555
##	KM (Spearman)	PAM (Euclidean)	HC (Euclidean)	KM (Euclidean)
##	0.5702522	0.6012224	0.6480550	0.6633613
##	PAM (MI)	KM (MI)		
##	0.8758883	0.9107256		

4.2 Kappa Statistics

The unadjusted kappa statistic is 0.5927477 and the weighted kappa statistic is 0.7717546 for the final meta consensus cluster.

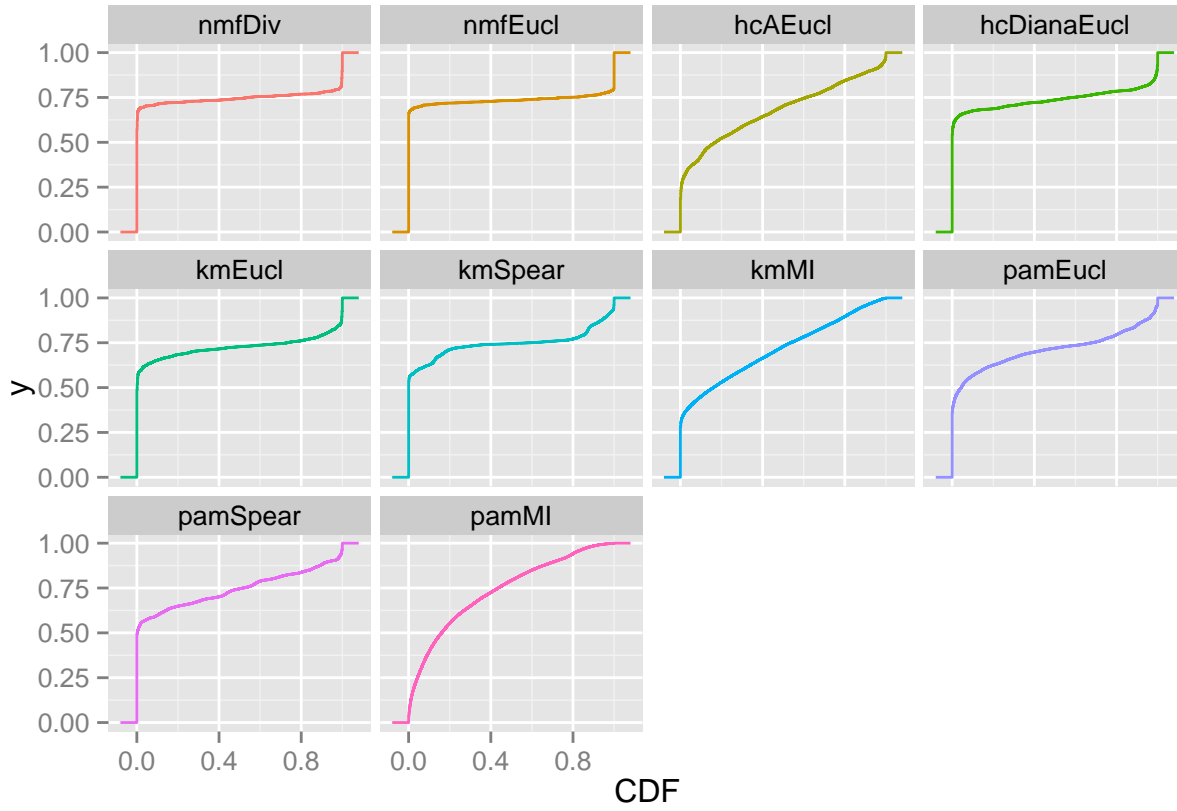
4.3 Adjusted Rand Index

The larger the better.

Algorithms	ARI
NMF (Divergence)	0.4799
NMF (Euclidean)	0.4435
PAM (Spearman)	0.427
HC (Diana)	0.4221
KM (Spearman)	0.4049
PAM (Euclidean)	0.3434
HC (Euclidean)	0.3275
KM (Euclidean)	0.2559
PAM (MI)	0.07465
KM (MI)	0.07369

5 Internal Evaluation

5.1 Proportion of Ambiguously Clustered Pairs (PAC)



The lower the better.

Algorithms	PAC
NMF (Euclidean)	0.1406
NMF (Divergence)	0.3267
HC (Diana)	0.3829
KM (Spearman)	0.425
KM (Euclidean)	0.4398
PAM (Spearman)	0.4748

Algorithms	PAC
PAM (Euclidean)	0.5984
KM (MI)	0.7167
HC (Euclidean)	0.7836
PAM (MI)	0.9638

The PAC for the meta consensus matrix is 0.7114804.

5.2 Davies-Bouldin Index

For DBI, the lower the better.

Algorithms	DBI
HC (Euclidean)	1.689
NMF (Euclidean)	1.702
NMF (Divergence)	1.702
PAM (Spearman)	1.72
PAM (Euclidean)	1.725
KM (Spearman)	1.725
HC (Diana)	1.731
KM (Euclidean)	1.741
PAM (MI)	1.935
KM (MI)	1.972

5.3 Dunn Index

For DI, the larger the better.

Algorithms	DI
HC (Euclidean)	1.04
PAM (Euclidean)	1.028
NMF (Euclidean)	1.003
NMF (Divergence)	0.999
KM (Spearman)	0.9926
HC (Diana)	0.9866
PAM (MI)	0.9821
PAM (Spearman)	0.9492
KM (MI)	0.9262
KM (Euclidean)	0.9238

5.4 Rousseeuw's Silhouette

For Rousseeuw's Silhouette internal cluster quality index (RS), the larger the better.

Algorithms	RS
HC (Euclidean)	0.1226
PAM (Euclidean)	0.1171

Algorithms	RS
NMF (Divergence)	0.1111
NMF (Euclidean)	0.1084
PAM (Spearman)	0.1069
KM (Spearman)	0.1059
HC (Diana)	0.09062
KM (Euclidean)	0.08707
PAM (MI)	-0.0061
KM (MI)	-0.007504

5.5 Silhouette Average Width

For SAW, the larger the better.

Algorithms	SAW
HC (Euclidean)	0.1226
PAM (Euclidean)	0.1171
NMF (Divergence)	0.1111
NMF (Euclidean)	0.1084
PAM (Spearman)	0.1069
KM (Spearman)	0.1059
HC (Diana)	0.09062
KM (Euclidean)	0.08707
PAM (MI)	-0.0061
KM (MI)	-0.007504

5.6 C-Index

For CI, the lower the better.

Algorithms	CI
KM (Euclidean)	0.2816
NMF (Divergence)	0.3023
PAM (Euclidean)	0.31
HC (Euclidean)	0.3129
HC (Diana)	0.3212
NMF (Euclidean)	0.3369
PAM (MI)	0.3376
PAM (Spearman)	0.3489
KM (MI)	0.3529
KM (Spearman)	0.356

5.7 Gamma Index

For BHI, the larger the better.

Algorithms	BHI
HC (Euclidean)	2.051

Algorithms	BHI
PAM (Euclidean)	1.754
HC (Diana)	1.746
PAM (Spearman)	1.715
KM (Euclidean)	1.676
NMF (Euclidean)	1.666
KM (Spearman)	1.627
NMF (Divergence)	1.617
KM (MI)	-3.142
PAM (MI)	-3.358

5.8 CH Index

For CHI, the larger the better.

Algorithms	CHI
NMF (Divergence)	80.36
NMF (Euclidean)	79.26
KM (Spearman)	78.71
PAM (Spearman)	77.19
PAM (Euclidean)	75.71
KM (Euclidean)	75.38
HC (Diana)	74.95
HC (Euclidean)	66.49
PAM (MI)	13.23
KM (MI)	7.322

5.9 Summary

Here is a summary of all the indices for each algorithm, in unsorted order.

Algorithms	PAC	DBI	DI	SAW	RS	CI	BHI	CHI
HC (Diana)	0.3829	1.731	0.9866	0.09062	0.09062	0.3212	1.746	74.95
HC (Euclidean)	0.7836	1.689	1.04	0.1226	0.1226	0.3129	2.051	66.49
KM (Euclidean)	0.4398	1.741	0.9238	0.08707	0.08707	0.2816	1.676	75.38
KM (MI)	0.7167	1.972	0.9262	-0.007504	-0.007504	0.3529	-3.142	7.322
KM (Spearman)	0.425	1.725	0.9926	0.1059	0.1059	0.356	1.627	78.71
NMF (Divergence)	0.3267	1.702	0.999	0.1111	0.1111	0.3023	1.617	80.36
NMF (Euclidean)	0.1406	1.702	1.003	0.1084	0.1084	0.3369	1.666	79.26
PAM (Euclidean)	0.5984	1.725	1.028	0.1171	0.1171	0.31	1.754	75.71
PAM (MI)	0.9638	1.935	0.9821	-0.0061	-0.0061	0.3376	-3.358	13.23
PAM (Spearman)	0.4748	1.72	0.9492	0.1069	0.1069	0.3489	1.715	77.19

6 Ranked Indices

The table below shows the ranking of algorithms for performance on a clustering index, for each index. There is an additional column that shows the proportion of indices where an algorithm was ranked **first** or

second.

Algorithms	PAC	DBI	DI	SAW	RS	CI	BHI	CHI	Top
HC (Euclidean)	9	1	1	1	1	4	1	8	62.5%
PAM (Euclidean)	7	5	2	2	2	3	2	5	50%
NMF (Divergence)	2	3	4	3	3	2	8	1	37.5%
NMF (Euclidean)	1	2	3	4	4	6	6	2	37.5%
KM (Euclidean)	5	8	10	8	8	1	5	6	12.5%
HC (Diana)	3	7	6	7	7	5	3	7	0%
KM (MI)	8	10	9	10	10	9	9	10	0%
KM (Spearman)	4	6	5	6	6	10	7	3	0%
PAM (MI)	10	9	7	9	9	7	10	9	0%
PAM (Spearman)	6	4	8	5	5	8	4	4	0%

If we were to conduct a meta consensus clustering, a weight for each algorithm needs to be assigned. One such way of doing so is using the inverse rank sums. We can sum the ranks for each algorithm, and assign higher weight for consistently high ranked methods (1st or 2nd) and vice versa. This is shown below:

Algorithms	Top	Sum	Weight
HC (Euclidean)	62.5%	26	14.69%
NMF (Divergence)	37.5%	26	14.69%
PAM (Euclidean)	50%	28	13.64%
NMF (Euclidean)	37.5%	28	13.64%
PAM (Spearman)	0%	44	8.68%
HC (Diana)	0%	45	8.49%
KM (Spearman)	0%	47	8.13%
KM (Euclidean)	12.5%	51	7.49%
PAM (MI)	0%	70	5.46%
KM (MI)	0%	75	5.09%

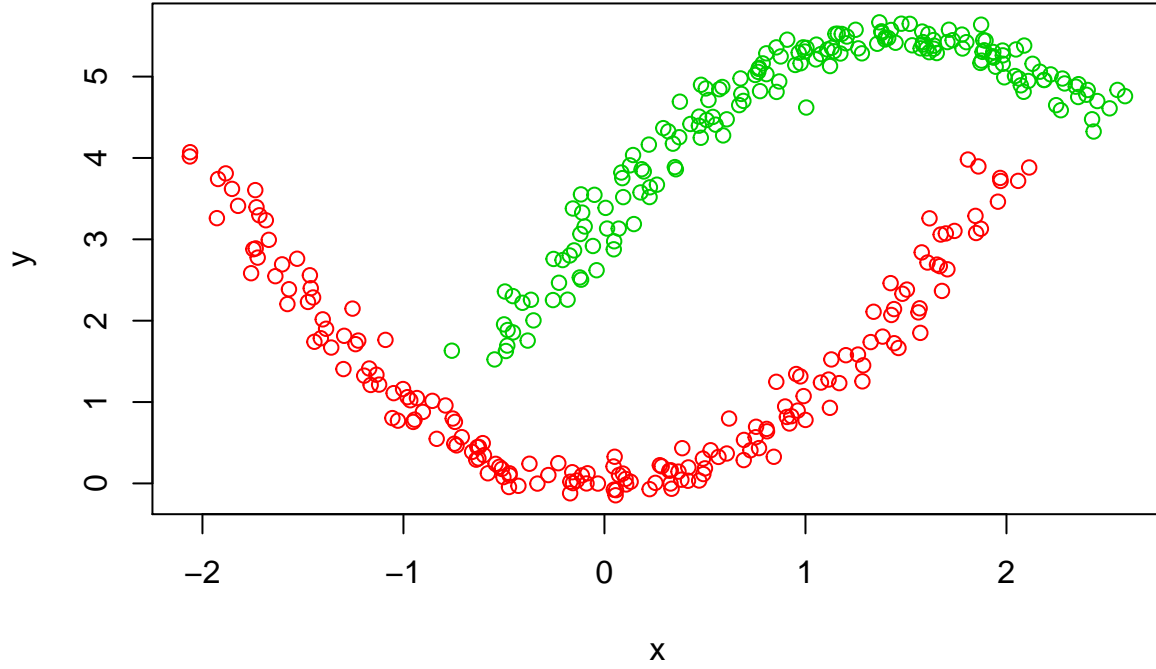
7 Simulations

To confirm the clustering results from using the TCGA dataset, we need to try out the algorithms on a few simulated datasets, designed to test the robustness of each method. The `clusterSim` package provides very good built-in examples to use.

7.1 Worms Dataset

The following two-cluster dataset is our first simulation:

Two clusters with atypical parabolic shapes (worms)



The ranked indices and weights tables are shown below:

Algorithms	PAC	DBI	DI	SAW	RS	CI	BHI	CHI	Top
KM (Euclidean)	6	2	1	1	1	1	1	1	87.5%
PAM (Euclidean)	9	1	4	2	2	2	2	2	75%
HC (Euclidean)	1	7	5	5	5	5	5	5	12.5%
HC (Diana)	1	7	5	5	5	5	5	5	12.5%
KM (Spearman)	1	7	5	5	5	5	5	5	12.5%
PAM (Spearman)	1	7	5	5	5	5	5	5	12.5%
NMF (Euclidean)	7	4	2	4	4	4	4	4	12.5%
KM (MI)	10	6	10	10	10	10	10	9	0%
PAM (MI)	5	5	9	9	9	9	9	10	0%
NMF (Divergence)	8	3	3	3	3	3	3	3	0%

Algorithms	Top	Sum	Weight
KM (Euclidean)	87.5%	14	22.9%
PAM (Euclidean)	75%	24	13.36%
NMF (Divergence)	0%	29	11.06%
NMF (Euclidean)	12.5%	33	9.72%
HC (Euclidean)	12.5%	38	8.44%
HC (Diana)	12.5%	38	8.44%

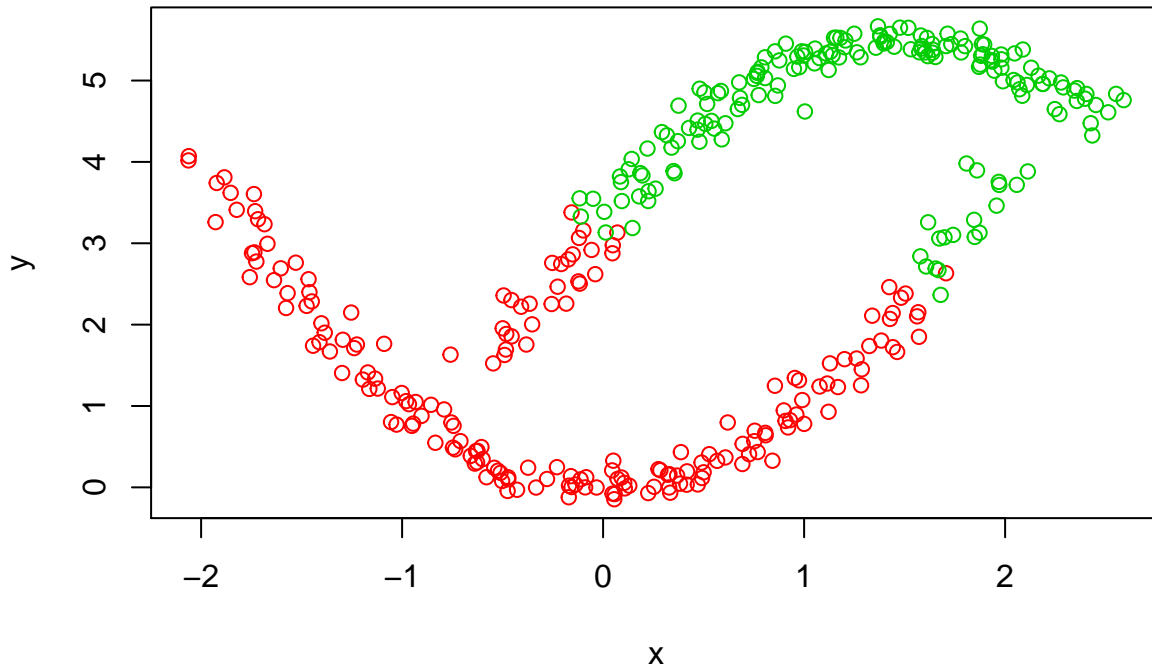
Algorithms	Top	Sum	Weight
KM (Spearman)	12.5%	38	8.44%
PAM (Spearman)	12.5%	38	8.44%
PAM (MI)	0%	65	4.93%
KM (MI)	0%	75	4.28%

Using hierarchical clustering with Wald’s method on the meta-consensus matrix across the ten algorithms, we can obtain meta-consensus classes. The following table shows how the different methods compare, based on the number of matches to the true class labels.

	Matches
KM (Euclidean)	311
PAM (Euclidean)	306
Meta	304
NMF (Divergence)	291
NMF (Euclidean)	287
PAM (MI)	181
KM (MI)	179
HC (Euclidean)	136
HC (Diana)	136
KM (Spearman)	136
PAM (Spearman)	136

Let’s visualize how the best algorithm, K-Means using euclidean distance, separates the clusters:

K–Means (euclidean) clustering of two atypical parabolic shapes (wor

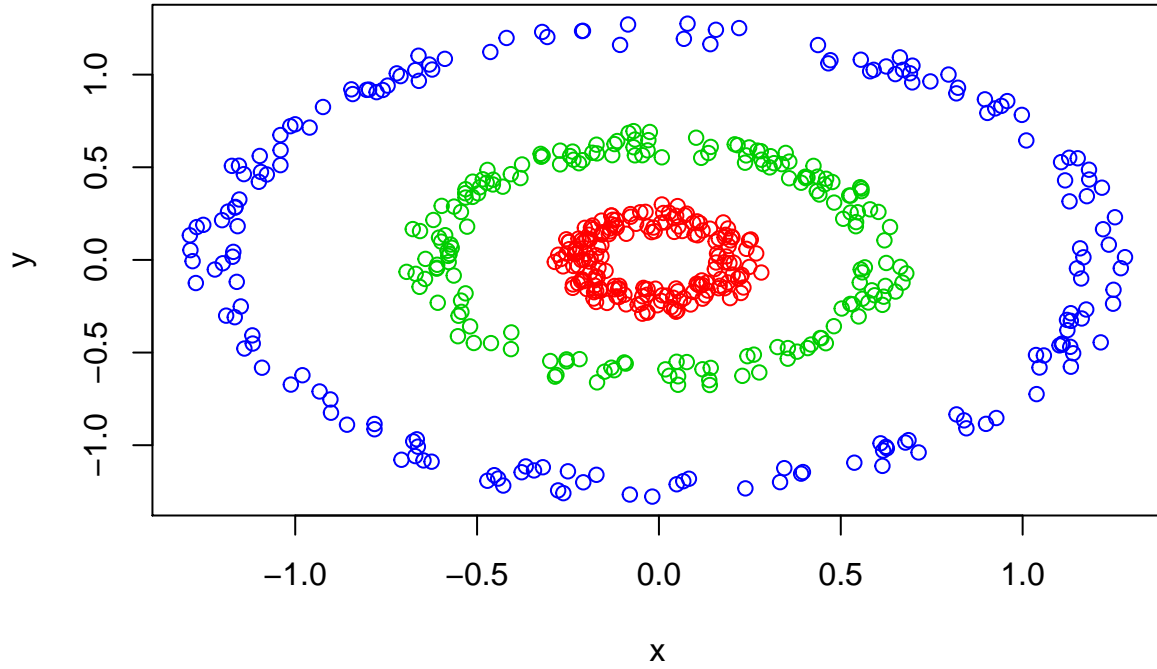


It appears that even the best algorithm cannot make a parabolic division in the feature space.

7.2 Rings Dataset

The following three-cluster dataset is our second simulation:

Three clusters with atypical ring shapes (circles)



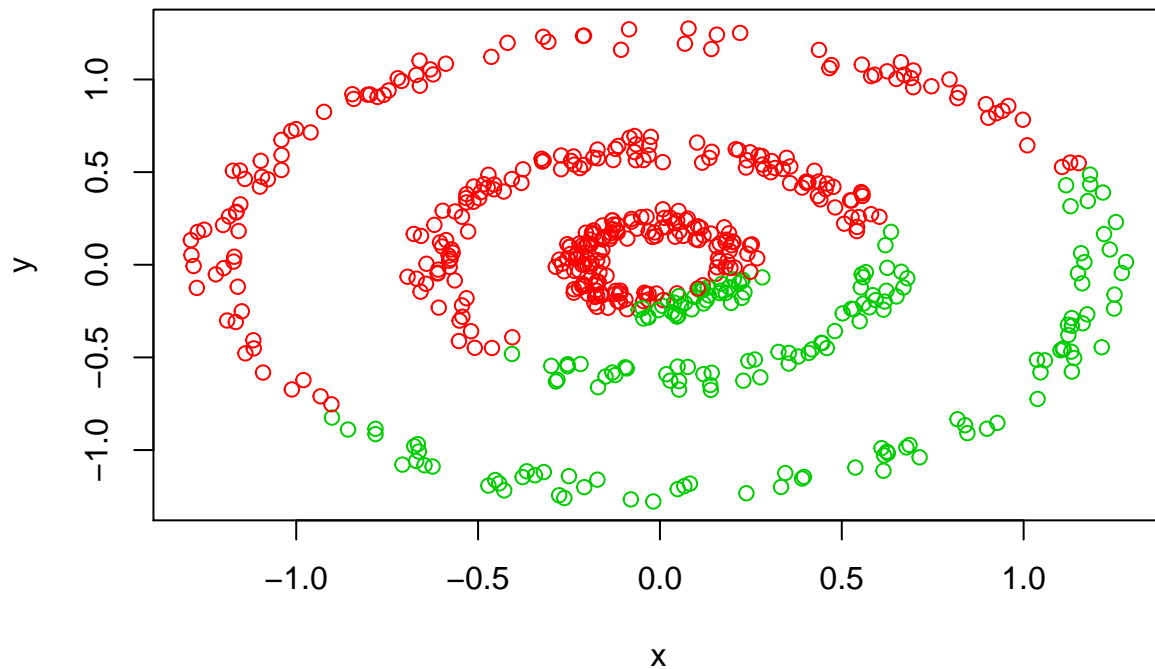
Algorithms	PAC	DBI	DI	SAW	RS	CI	BHI	CHI	Top
NMF (Euclidean)	6	2	1	2	2	1	8	1	75%
KM (Euclidean)	7	9	9	4	4	2	2	9	25%
KM (MI)	8	1	8	9	9	9	1	8	25%
PAM (Euclidean)	9	4	3	1	1	7	3	7	25%
NMF (Divergence)	5	3	2	3	3	8	9	2	25%
HC (Euclidean)	1	5	4	5	5	3	4	3	12.5%
HC (Diana)	1	5	4	5	5	3	4	3	12.5%
KM (Spearman)	1	5	4	5	5	3	4	3	12.5%
PAM (Spearman)	1	5	4	5	5	3	4	3	12.5%

Algorithms	Top	Sum	Weight
NMF (Euclidean)	75%	23	15.84%
HC (Euclidean)	12.5%	30	12.14%
HC (Diana)	12.5%	30	12.14%
KM (Spearman)	12.5%	30	12.14%
PAM (Spearman)	12.5%	30	12.14%
PAM (Euclidean)	25%	35	10.41%
NMF (Divergence)	25%	35	10.41%

Algorithms	Top	Sum	Weight
KM (Euclidean)	25%	46	7.92%
KM (MI)	25%	53	6.87%

	Matches
PAM (Euclidean)	201
KM (MI)	181
KM (Euclidean)	180
NMF (Euclidean)	180
HC (Euclidean)	169
HC (Diana)	169
KM (Spearman)	169
PAM (Spearman)	169
NMF (Divergence)	161
Meta	146

PAM (euclidean) clustering of three circles (rings)



Again, we are only able to make linear separations.

8 References

- Monti, S., Tamayo, P., Mesirov, J., & Golub, T. (2003). "Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data." *Machine learning*, 52 (1-2), 91-118.
- Conrad, J. G., Al-Kofahi, K., Zhao, Y., & Karypis, G. (2005, June). "Effective document clustering for large heterogeneous law firm collections." In *Proceedings of the 10th international conference on Artificial intelligence and law* (pp. 177-187). ACM.

- Rand, W. M. (1971). "Objective criteria for the evaluation of clustering methods." *Journal of the American Statistical Association*, 66 (336), 846-850.
- Hubert, L., & Arabie, P. (1985). "Comparing partitions." *Journal of Classification*, 2 (1), 193-218.
- Vinh, N. X., Epps, J., & Bailey, J. (2009, June). "Information theoretic measures for clusterings comparison: is a correction for chance necessary?." In *Proceedings of the 26th Annual International Conference on Machine Learning* (pp. 1073-1080). ACM.
- Senbabaoglu, Y., Michailidis, G., & Li, J. Z. (2014). "Critical limitations of consensus clustering in class discovery." *Scientific reports*, 4.
- Davies, David L.; Bouldin, Donald W. (1979). "A Cluster Separation Measure." *IEEE Transactions on Pattern Analysis and Machine Intelligence*. PAMI-1 (2): 224-227.
- Dunn, J. C. (1973). "A Fuzzy Relative of the ISODATA Process and Its Use in Detecting Compact Well-Separated Clusters." *Journal of Cybernetics* 3 (3): 32-57.
- Rousseeuw, P. J. (1987). "Silhouettes: a graphical aid to the interpretation and validation of cluster analysis." *Journal of Computational and Applied Mathematics*, 20, 53-65.
- Hubert, L. J., & Levin, J. R. (1976). "A general statistical framework for assessing categorical clustering in free recall." *Psychological bulletin*, 83 (6), 1072-1080.
- Baker, F. B., & Hubert, L. J. (1975). "Measuring the power of hierarchical cluster analysis." *Journal of the American Statistical Association*, 70 (349), 31-38.
- Calinski, T., & Harabasz, J. (1974). "A dendrite method for cluster analysis." *Communications in Statistics-theory and Methods*, 3 (1), 1-27.