



Clustering Leukemia Patients

Katelyn Gao¹, Heather Hardeman², Edward Lim³, Cris Potter⁴, Carl Meyer⁵, Ralph Abbey⁵



Introduction

- **Clustering** partitions a set of observations into groups where similar observations are grouped together.
- Our Clustering Methods
 - **k-means**
 - $\arg \min_S \sum_{i=1}^k \sum_{x_j \in S_i} \|x_j - \mu_i\|^2$
 - Initialization \rightarrow Assignment and Update \rightarrow Termination
 - **Nonnegative Matrix Factorization (NMF)**
 - Factor a nonnegative data matrix into two nonnegative matrices

$$\min \|V_{m \times n} - W_{m \times k} H_{k \times n}\|_F^2$$

Update Steps

$$H_{ij} \leftarrow H_{ij} \frac{(W^T V)_{ij}}{(W^T W H)_{ij}} \quad W_{ij} \leftarrow W_{ij} \frac{(V H^T)_{ij}}{(W H H^T)_{ij}}$$

Motivation

Determine the number of clusters present in a data set and use that information to cluster Leukemia patients

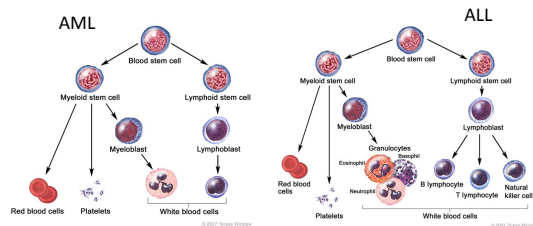
Data from the Broad Institute

Types of **Leukemia**

- Acute Lymphoblastic Leukemia B cells (ALL - B)
- Acute Lymphoblastic Leukemia T-cells (ALL - T)
- Acute Myeloid Leukemia (AML)

5000 genes x 38 leukemia patients' microarray data

- 1-19: ALL - B
- 20-27: ALL - T
- 28-38: AML



Methodology

Steps: NMF and K-means \rightarrow Consensus Matrix \rightarrow Sinkhorn Knopp \rightarrow Eigen decomposition

Consensus Matrix

- For Leukemia data set $A = 38 \times 38$ matrix
- $A: a_{ij} = 1$ if i and j are in the same cluster
- $a_{ij} = 0$ if not
- Symmetric
- Consensus matrix: for n runs of a clustering method, $C = \frac{1}{n} \sum_{i=1}^n A_i$

Sinkhorn Knopp Algorithm

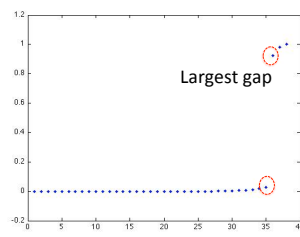
- Create a matrix S by adjusting column sums and row sums of the consensus matrix to 1
- Alternate scaling rows by their row sums and the columns by their column sums
- Preserves the symmetry of a matrix and places the eigenvalues in the interval $[0, 1]$
- Preserves the block structure of a consensus matrix

Results

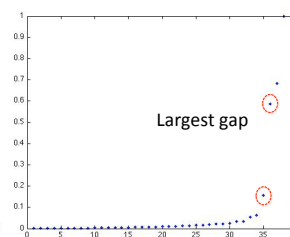
Eigenvalues

We used the location of the largest gap in the eigenvalues of the matrix S to determine the number of clusters k . Here the largest gap occurs between the third and fourth, implying $k = 3$

50 runs of NMF $k = 3$

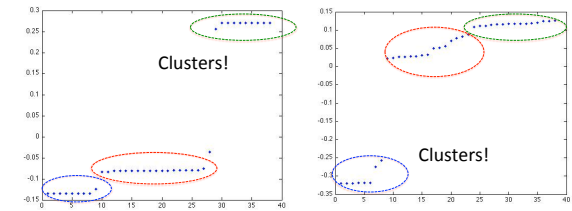


50 runs of K-means $k = 2, 3, 4, 5$



Eigenvectors

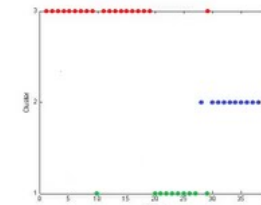
The eigenvectors may reveal information about which cluster each observation belongs.



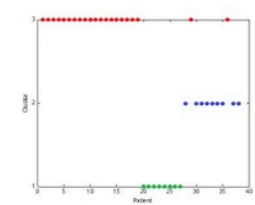
Conclusions

- The location of the largest gap in the eigenvalues reveals the number of clusters k .
- Then we may re-cluster the data based on that information.

NMF with $k = 3$



k-means with $k = 3$



- Our results corroborate with those of the Broad Institute of Harvard and MIT.
- Information from the Broad Institute suggests that the "mis-clustering" of patients 10 and 29 is due to misdiagnoses.

Acknowledgements

Dr. Carl D. Meyer, *advisor*
Ralph Abbey, *graduate assistant*
Broad Institute of Harvard and MIT