

M4: Unsupervised ML methods

- Immediate applications of (spectral) linear algebra

Manikandan Narayanan

Week 8 (Sep 15-, 2025)

PRML Jul-Nov 2025 (Grads Section)

Acknowledgment of Sources

- Slides based on content from related
 - Courses:
 - IITM – Profs. Arun/Harish[HR]/Chandra[CC]/Prashanth’s PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi’s “Intro to ML” slides – cited (e.g., [HR]/[HG]) in the bottom right of a slide.
 - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.
 - Books:
 - PRML by **Bishop**. (content, figures, slides, etc.) – cited as [CMB]
 - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
 - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
 - Foundations of ML by Mohri, Rostamizadeh, and Talwalkar (content, figures, slides by Mohri, etc.). – [MRT]
 - Information Theory, Inference and Learning Algorithms by MacKay – [DJM]

Outline for Module M4

- **M4.0 Context/Introduction**
- M4.1 Spectral clustering algorithm
 - M4.1.1 From data to similarity graph
 - M4.1.2 From graph connectivity to clusters
- M4.2 Conclusion & next steps

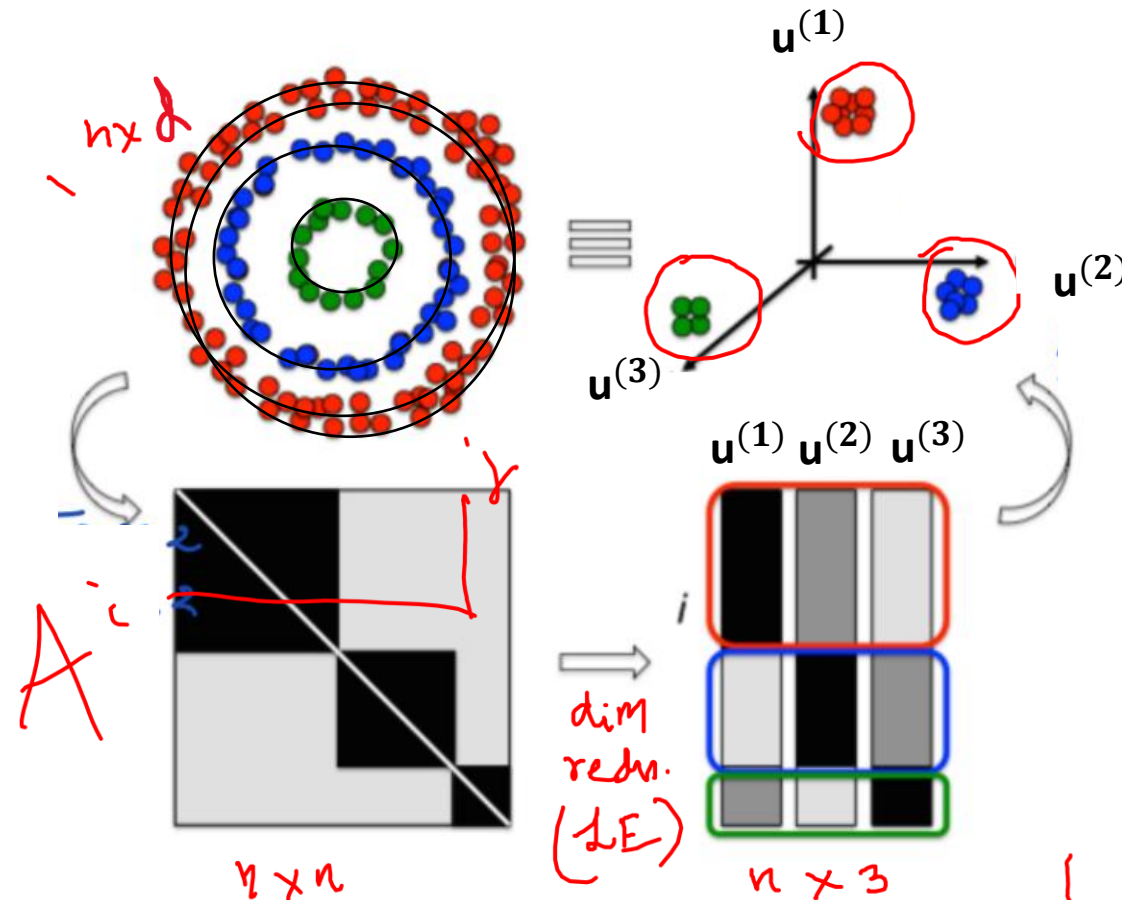
Context/Intro.: Two unsupervised ML problems

- Unsupervised ML: Recognize patterns in the dataset (set of n data points in \mathbb{R}^d) without any labels on the data points.
- Dimensionality reduction:
 - Transforms data points from high to low dimensions without much loss of “information”, assuming the data points lie “effectively in/close to” a low-dim. manifold of the original space.
 - Many approaches possible: **PCA**, t-SNE, UMAP, **Laplacian Eigenmaps (spectral)**, etc.
- Clustering:
 - Grouping n objects into k clusters based on their similarity
 - Again, many approaches possible: **k-means**, hierarchical, **spectral**, etc.

Outline for Module M4

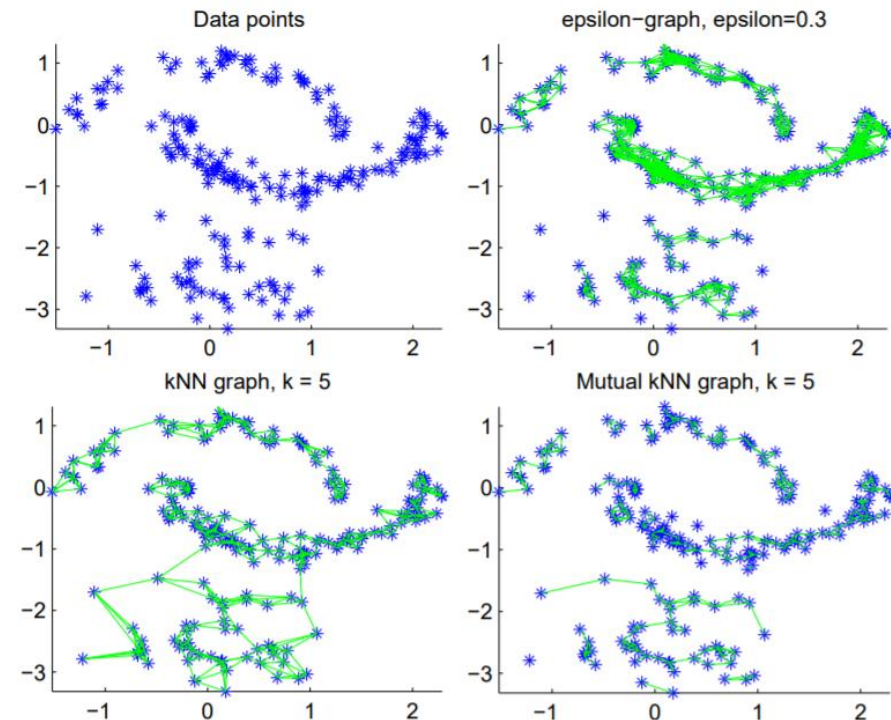
- M4.0 Context/Introduction
- **M4.1 Spectral clustering algorithm**
 - **M4.1.1 From data to similarity graph**
 - M4.1.2 From graph connectivity to clusters
- M4.2 Conclusion & next steps

Spectral clustering == dim. redn. + kmeans



Step 1: From $(n \times d)$ data matrix to
 $(n \times n)$ similarity-graph's adj. matrix!

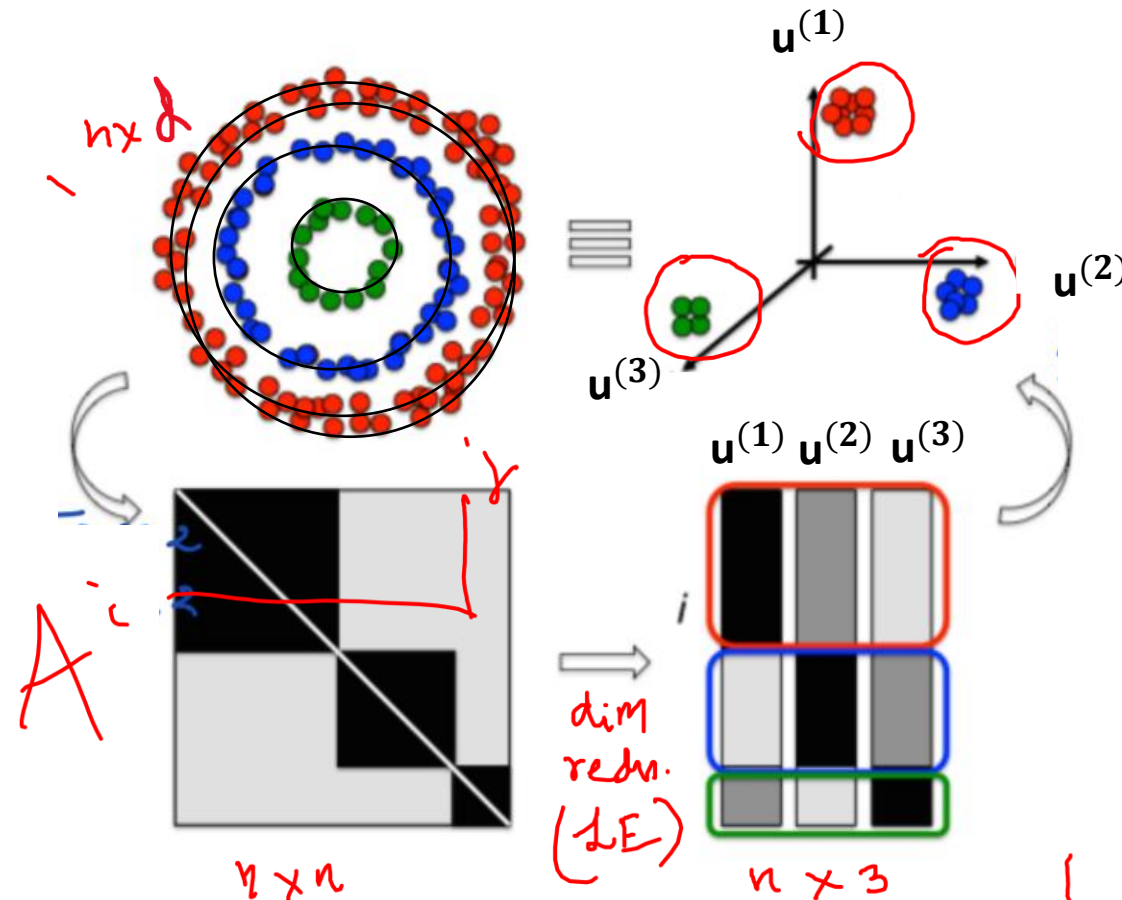
- Similarity function: $s(x_i, x_j) = \exp(-\|x_i - x_j\|^2 / (2\sigma^2))$
- Similarity graph: complete (all-pairs) graph, or sparse graphs



Outline for Module M4

- M4.0 Context/Introduction
- **M4.1 Spectral clustering algorithm**
 - M4.1.1 From data to similarity graph
 - **M4.1.2 From graph connectivity to clusters**
- M4.2 Conclusion & next steps

Spectral clustering == dim. redn. + kmeans



Step 2. Spectral clustering – problem defn. using graph connectivity (i.e., adjacency matrix A)

Input:

$A \in \mathbb{R}^{n \times n}$: A_{ij} = Similarity between datapoint i & j
 ≥ 0

Output:

$U \in \mathbb{R}^n$: A score assignment to all nodes.

e.g. $cl. 2 [0 \ 0 \ 0 \ 1 \ 1 \ 1]$
 $cl. 1 [1 \ 1 \ 1 \ 0 \ 0 \ 0]$

~Obj. fn:

$$L(u) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{ij} (u_i - u_j)^2$$

e.g.

$$A = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$

clust1 clust2
 $\checkmark [a \ a \ a \ b \ b \ b]$
 $\times [a \ a \ a \ a \ a \ a]$

When is $L(u)$ small?

~Obj. fn.
 (constrained):

$$\min_u L(u) \text{ s.t. } \underline{u^T u = 1}$$

$$D = \begin{bmatrix} 3 & 0 \\ 0 & \dots \end{bmatrix}$$

Laplacian Matrix ($L=D-A$)

$$L(u) = \sum_{i=1}^n \sum_{j=1}^n u_i u_j (-A_{ij}) + \sum_{i=1}^n u_i^2 \sum_{j=1}^n A_{ij}$$

$$= u^T (D - A) u$$

$$\text{where } D_{ii} = \sum_j A_{ij}$$

$$\text{Laplacian matrix } L = D - A$$

Blank space: Example Graph Laplacian

e.g.

$$A = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}_{n \times n}$$

sim. graph.

$$D = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{bmatrix}_{n \times n}$$


$$L = D - A = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{bmatrix}$$

Background on Constrained Optimization (i): Lagrange multiplier (single constraint)

Problem: opt (min or max) $f(\mathbf{x})$ $(\mathbf{x} \in \mathbb{R}^n, f: \mathbb{R}^n \rightarrow \mathbb{R})$
s.t. constraint $g(\mathbf{x}) = c$

FONC (First-order necessary conditions): Instead of setting gradient of $f(\mathbf{x})$ to zero as in unconstrained optima, set gradient of Lagrangian function $\mathcal{L}(\mathbf{x}, \lambda)$ to zero for constrained optima!

$$\nabla \mathcal{L}(\mathbf{x}, \lambda) = 0, \quad \text{where } \mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda(c - g(\mathbf{x})); \lambda \in \mathbb{R}.$$

Opt. strategy: Equivalently, we've these $(n+1)$ equations over $(n+1)$ unknowns to solve, and inspect each solution to find the constrained maxima or minima.

$$\begin{aligned} \nabla f(\mathbf{x}) &= \lambda \nabla g(\mathbf{x}) && \text{(n equations)} \\ g(\mathbf{x}) &= c && \text{(1 equation)} \end{aligned}$$

Fineprint: For Lagrange multiplier to work, f, g must be smooth (continuously differentiable) and constraint qualification $\nabla g(\mathbf{x}^*) \neq 0$ must hold. Otherwise, solution may not exist for the above $(n+1)$ equations!

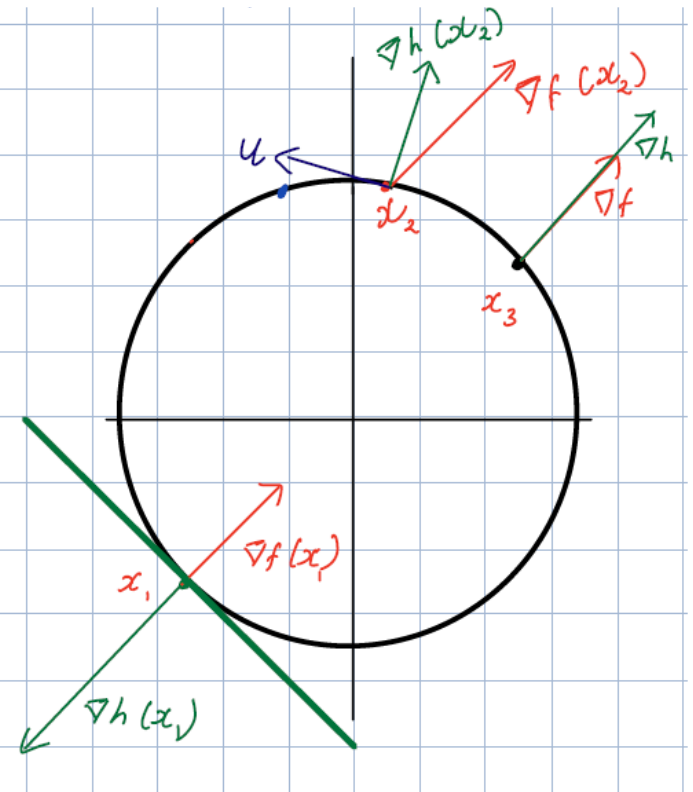
Background on Constrained Optimization (i): Example

e.g.

$$\min x+y$$

s.t.

$$x^2 + y^2 - 1 = 0$$



How do you minimize $L(\mathbf{u}^{(1)})$ subject to unit norm constraint?

Minimize $\mathbf{u}^{(1)} \in \mathbb{R}^n \quad \mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} \quad \text{s.t.} \quad \mathbf{u}^{(1)T} \mathbf{u}^{(1)} = 1$

Set gradient of Lagrangian: $\mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} + \lambda_1 (1 - \mathbf{u}^{(1)T} \mathbf{u}^{(1)})$ to zero to

get these eqns: $\mathbf{L} \mathbf{u}^{(1)} = \lambda_1 \mathbf{u}^{(1)}$
 $\mathbf{u}^{(1)T} \mathbf{u}^{(1)} = 1$

$$\nabla_x (x^T A x) = (A + A^T) x$$
$$\nabla_x (x^T x) = 2x$$

Thus, any solution $\mathbf{u}^{(1)}$ is an unit-norm eigen vector of \mathbf{L} , with function value being:

$$\mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} = \mathbf{u}^{(1)T} \lambda_1 \mathbf{u}^{(1)} = \lambda_1 \mathbf{u}^{(1)T} \mathbf{u}^{(1)} = \lambda_1$$

To minimize above function value, choose the solution (eigen vector) with the smallest eigen value!

Given that optimal $u^{(1)}$ is an eigen vector of L , what do we know about the spectra of L ?

Let's look at the spectral properties of the Laplacian matrix L ...

Properties of Laplacian $L=D-A$: Initial obsns.

L (real). Symm.

L psd $\Rightarrow \lambda_i \geq 0$

Is $\lambda_i = 0$?

If $u = [1, 1, \dots, 1]^T$, then $L(u) = 0$ ($\lambda_i = 0$)

Properties of Laplacian: Toy Graph Example

e.g.

$$A = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix}_{n \times n}$$

sim. graph.

$$D = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{bmatrix}_{n \times n}$$


$$L \vec{x} = \lambda \vec{x}$$

$$L = D - A = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & -1 & 2 \end{bmatrix}$$

orthog.

$$\begin{aligned} \vec{u}^{(1)} &= [1 \ 1 \ 1 \ 0 \ 0 \ 0] & \lambda &= 0 \\ \vec{u}^{(2)} &= [0 \ 0 \ 0 \ 1 \ 1 \ 1] & \lambda &= 0 \\ \vec{u} &= [1 \ 1 \ 1 \ 1 \ 1 \ 1] & \lambda &= 0 \end{aligned}$$

$$\begin{aligned} L \vec{u}^{(1)} &= \underline{0} \cdot \vec{u}^{(1)} \\ L \vec{u}^{(2)} &= \underline{0} \cdot \vec{u}^{(2)} \\ &\vdots \end{aligned}$$

Properties of Laplacian $L=D-A$: formally

Proposition 1 (Properties of L) *The matrix L satisfies the following properties:*

1. For every vector $u \in \mathbb{R}^n$ we have

$$u^T L u = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{ij} (u_i - u_j)^2$$

2. L is symmetric and positive semi-definite.

3. The smallest eigenvalue of L is 0, the corresponding eigenvector is the constant one vector $\mathbf{1}$.

4. L has n non-negative, real-valued eigenvalues $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$.

Proposition 2 (Number of connected components and the spectrum of L) *Let G be an undirected graph with non-negative weights. Then the multiplicity k of the eigenvalue 0 of L equals the number of connected components A_1, \dots, A_k in the graph. The eigenspace of eigenvalue 0 is spanned by the indicator vectors $\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_k}$ of those components.*

(Prove for $k=1$)
(Then use spectrum of block diag. matrices)

Going from one ($\mathbf{u}^{(1)}$) to two scores per node!

- What is the objective function and constraints with
 - one score vector?
 - two score vectors $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}$?
 - (find two score vectors that minimize $\mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} + \mathbf{u}^{(2)T} \mathbf{L} \mathbf{u}^{(2)}$ s.t. orthonormality)
 - three score vectors?
 - ...
 - M score vectors?

Background on Constrained Optimization (ii): Lagrange multiplier (multiple constraints)

Problem: opt (min or max) $f(\mathbf{x})$ $(\mathbf{x} \in \mathbb{R}^n, f: \mathbb{R}^n \rightarrow \mathbb{R})$
s.t. constraints $g_1(\mathbf{x}) = 0, \dots, g_m(\mathbf{x}) = 0$ ($m < n$)

FONC (First-order necessary conditions): same: set gradient of Lagrangian function $\mathcal{L}(\mathbf{x}, \lambda)$ to zero for constrained optima!
 $\nabla \mathcal{L}(\mathbf{x}, \lambda) = 0$, where $\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x}), \lambda_i \in \mathbb{R}$.

Opt. strategy: Equivalently, we've these $(n+m)$ equations over $(n+m)$ unknowns to solve, and inspect each solution to find the constrained maxima or minima.

$$\begin{aligned} \nabla f(\mathbf{x}) &= - \sum_i \lambda_i \nabla g_i(\mathbf{x}) && \text{(n equations)} \\ g_i(\mathbf{x}) &= 0 && \text{(m equations)} \end{aligned}$$

Fineprint: For Lagrange multiplier to work, f, g_i must be smooth (continuously differentiable) and constraint qualification " $\nabla g(x^*) \neq 0$ " becomes "the constraint gradient vectors $\{\nabla g_i(x^*)\}_i$ should be linearly indept.". (Otherwise, solution may not exist again).

How do you minimize $L(\mathbf{u}^{(1)}) + L(\mathbf{u}^{(2)})$ subject to orthonormality constraints?

$$\begin{aligned} \text{Minimize }_{\mathbf{u}^{(1)}, \mathbf{u}^{(2)} \in \mathbb{R}^n} f(\mathbf{u}^{(1)}, \mathbf{u}^{(2)}) &:= \mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} + \mathbf{u}^{(2)T} \mathbf{L} \mathbf{u}^{(2)} \\ \text{s.t. } \mathbf{u}^{(1)T} \mathbf{u}^{(1)} &= 1, \mathbf{u}^{(2)T} \mathbf{u}^{(2)} = 1, \text{ and} \\ \mathbf{u}^{(1)T} \mathbf{u}^{(2)} &= 0 \text{ (equivalently } -2 \mathbf{u}^{(1)T} \mathbf{u}^{(2)} = 0 \text{ for convenience)} \end{aligned}$$

Set gradient of Lagrangian: $\mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} + \mathbf{u}^{(2)T} \mathbf{L} \mathbf{u}^{(2)} + \lambda_1(1 - \mathbf{u}^{(1)T} \mathbf{u}^{(1)}) + \lambda_2(1 - \mathbf{u}^{(2)T} \mathbf{u}^{(2)}) + \lambda_{12}(-2 \mathbf{u}^{(1)T} \mathbf{u}^{(2)})$ to zero to

get these eqns:

$$\partial L / \partial \mathbf{u}^{(1)} = 2\mathbf{L}\mathbf{u}^{(1)} - 2\lambda_1 \mathbf{u}^{(1)} - 2\lambda_{12} \mathbf{u}^{(2)} = 0 \quad \text{-(1)}$$

$$\partial L / \partial \mathbf{u}^{(2)} = 2\mathbf{L}\mathbf{u}^{(2)} - 2\lambda_2 \mathbf{u}^{(2)} - 2\lambda_{12} \mathbf{u}^{(1)} = 0 \quad \text{-(2)}$$

$$\mathbf{u}^{(1)T} \mathbf{u}^{(1)} = 1, \mathbf{u}^{(2)T} \mathbf{u}^{(2)} = 1, \mathbf{u}^{(1)T} \mathbf{u}^{(2)} = 0 \quad \text{-(3)}$$

One solution is given by: $\lambda_{12} = 0$, and $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}$ being any two orthonormal eigen vectors of \mathbf{L} , with function value f being the sum of their corresp. eigen values: $\lambda_1 + \lambda_2$.

Thus to minimize f , choose the two orthonormal eigen vectors of \mathbf{L} with the smallest two eigen values (incl. multiplicity)!

Final check: Qn. What about other solutions where $\lambda_{12} \neq 0$?

Ans. They can be converted to a solution with $\lambda_{12} = 0$ without changing the function value!

$$\text{Let } \mathbf{U}_{n \times 2} = [\mathbf{u}^{(1)} \ \mathbf{u}^{(2)}] \quad \text{and} \quad H = \begin{bmatrix} \lambda_1 & \lambda_{12} \\ \lambda_{12} & \lambda_2 \end{bmatrix}.$$

Then, Eqns. (1)-(3) can be rewritten in matrix notation as:

$$\mathbf{L}\mathbf{U} = \mathbf{U}H \quad \text{-(4)}$$

$$\mathbf{U}^T \mathbf{U} = I \text{ (identity matrix)} \quad \text{-(5)}$$

$$\text{, and the function value can be written as a matrix trace: } f(\mathbf{U}) = \mathbf{u}^{(1)T} \mathbf{L} \mathbf{u}^{(1)} + \mathbf{u}^{(2)T} \mathbf{L} \mathbf{u}^{(2)} = \text{Tr}(\mathbf{U}^T \mathbf{L} \mathbf{U}) \quad \text{-(6)}$$

Fact: Any solution $(\mathbf{U} = \hat{\mathbf{U}}, H = \hat{H})$ to eqns. (4)-(5) with \hat{H} a non-diagonal matrix (i.e., $\lambda_{12} \neq 0$) can be converted to a solution $(\mathbf{U} = \hat{\mathbf{U}}Q, H = \Lambda)$ with Λ a diagonal matrix (i.e., $\lambda_{12} = 0$) and the same function value (i.e., $f(\hat{\mathbf{U}}Q) = f(\hat{\mathbf{U}})$). Here, Q, Λ can be computed from orthogonal diagonalization of the symmetric matrix \hat{H} as $\hat{H} = Q\Lambda Q^T$.

Proof of Fact: Given that $(\hat{\mathbf{U}}, \hat{H})$ satisfies Eqns. (4)-(5), we first show that $(\hat{\mathbf{U}}Q, \Lambda)$ also satisfies Eqns. (4)-(5) below.

- $\mathbf{L}\hat{\mathbf{U}} = \hat{\mathbf{U}}\hat{H} \Rightarrow \mathbf{L}\hat{\mathbf{U}} = \hat{\mathbf{U}}Q\Lambda Q^T \Rightarrow \mathbf{L}(\hat{\mathbf{U}}Q) = (\hat{\mathbf{U}}Q)\Lambda$
- $(\hat{\mathbf{U}}Q)^T \hat{\mathbf{U}}Q = Q^T \hat{\mathbf{U}}^T \hat{\mathbf{U}}Q = Q^T Q = I$

Now, the function value (6) of any solution is: $f(\mathbf{U}) = \text{Tr}(\mathbf{U}^T \mathbf{L} \mathbf{U}) = \text{Tr}(\mathbf{U}^T \mathbf{U} H)$ (by Eqn (4)) = $\text{Tr}(H)$ (by Eqn (5)). Therefore,

- $f(\hat{\mathbf{U}}Q) = \text{Tr}(\Lambda)$, and $f(\hat{\mathbf{U}}) = \text{Tr}(\hat{H})$. But, $\text{Tr}(\Lambda) = \text{Tr}(\hat{H})$ because the sum of eig.vals. of any matrix is its trace.

Going from two $(\mathbf{u}^{(1)}, \mathbf{u}^{(2)})$ to multiple scores per node!

- What is the objective function and constraints with
 - one score vector?
 - two score vectors?
 - three score vectors?
 - ...
 - M score vectors $\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \dots, \mathbf{u}^{(M)}$?
 - (find M score vectors that minimize $\sum_k \mathbf{u}^{(k)T} \mathbf{L} \mathbf{u}^{(k)}$ s.t. orthonormality)

How do you minimize $\sum_k L(\mathbf{u}^{(k)})$ subject to orthonormality constraints? Outline of answer (without proof)

Let $\mathbf{U}_{n \times M} = [\mathbf{u}^{(1)} \ \mathbf{u}^{(2)} \ \dots \ \mathbf{u}^{(M)}]$, and Lagrange multipliers $H_{M \times M} = \{\lambda_{ij}\}_{i,j=1,\dots,M}$.

Problem is to Minimize $f(\mathbf{U}) = \sum_{k=1}^M \mathbf{u}^{(k)T} \mathbf{L} \mathbf{u}^{(k)} = \text{Tr}(\mathbf{U}^T \mathbf{L} \mathbf{U})$
s.t. $\mathbf{U}^T \mathbf{U} = \mathbf{I}$

Setting gradient of Lagrangian $\text{Tr}(\mathbf{U}^T \mathbf{L} \mathbf{U}) + \text{Tr}(\mathbf{H}(\mathbf{I} - \mathbf{U}^T \mathbf{U}))$ to zero yields:

$$\mathbf{L} \mathbf{U} = \mathbf{U} \mathbf{H} \quad \text{-(7)}$$

$$\mathbf{U}^T \mathbf{U} = \mathbf{I} \text{ (identity matrix)} \quad \text{-(8)}$$

Eqns. (7)-(8) can be solved and constrained minima found as in the $M=2$ case! Specifically,

- H can be taken to be a symmetric matrix, and further assumed wlog (without loss of generality) to be diagonal.
- Then, any solution should be M orthonormal eigen vectors of L .
- Minimizer found by taking the eigen vectors corresp. to the M smallest eigen values.

[related to Exercise 12.2 in [CMB]]

In other words, solution using $M=(a+1)$ th smallest eigen values (and corresp. orthonormal eigen vectors)

Let $u^0, u^1 \dots u^a$ be the eigen vectors corresponding to $\lambda_0, \lambda_1, \dots, \lambda_a$ the $a+1$ smallest eigen-values of L .

Then

$$\sum_{i,j} A_{ij} (u_i^k - u_j^k)^2 = u^{kT} L u^k = \lambda_k$$

Thus if λ_k is small, then the score assigned by u^k is such that $\text{Large } A_{ij} \Rightarrow \text{Small } |u_i^k - u_j^k|$



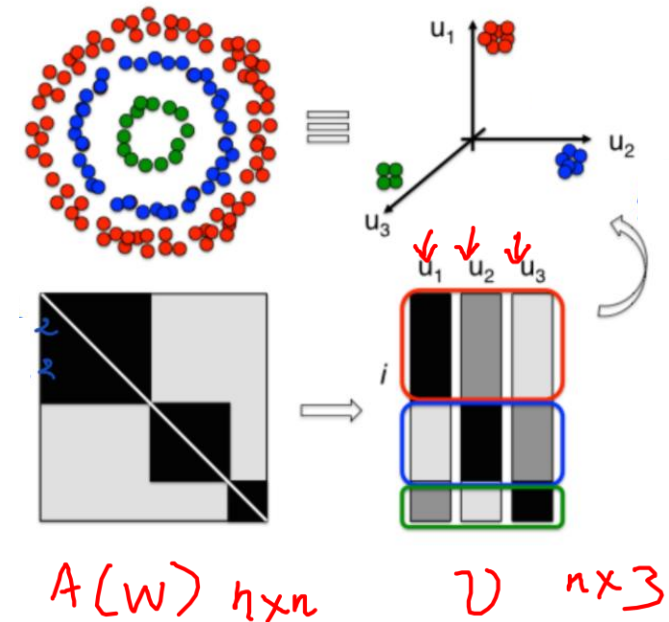
Spectral clustering algorithm:

Unnormalized spectral clustering

Input: Similarity matrix $S \in \mathbb{R}^{n \times n}$, number k of clusters to construct.

- Construct a similarity graph by one of the ways described in Section 2. Let W be its weighted adjacency matrix.
- Compute the unnormalized Laplacian L .
- Compute the first k eigenvectors u_1, \dots, u_k of L .
- Let $U \in \mathbb{R}^{n \times k}$ be the matrix containing the vectors u_1, \dots, u_k as columns.
- For $i = 1, \dots, n$, let $y_i \in \mathbb{R}^k$ be the vector corresponding to the i -th row of U .
- Cluster the points $(y_i)_{i=1, \dots, n}$ in \mathbb{R}^k with the k -means algorithm into clusters C_1, \dots, C_k .

Output: Clusters A_1, \dots, A_k with $A_i = \{j \mid y_j \in C_i\}$.



Conclusion and Looking ahead

- Linear Algebra has some immediate applications in Unsupervised ML.
 - We have looked at spectral-related applications (for both dim. redn. and clustering), by inspecting the eigen-values/vectors (spectra) of the Laplacian matrix.
 - Spectral clustering code on toy examples:
https://colab.research.google.com/drive/1b3AEjJ2_QseWJrGBwWb1x_08OVBXwZS3?usp=sharing .
- Next steps:
 - Further applications of Linear Algebra:
 - Dim. redn (contd. appn. of spectra of another matrix) – PCA
 - Supervised ML (appn. of projection in linear algebra) – Linear regression

Thank you!

- Appendix follows

Defn. of convexity of a cluster

- **Observation:**

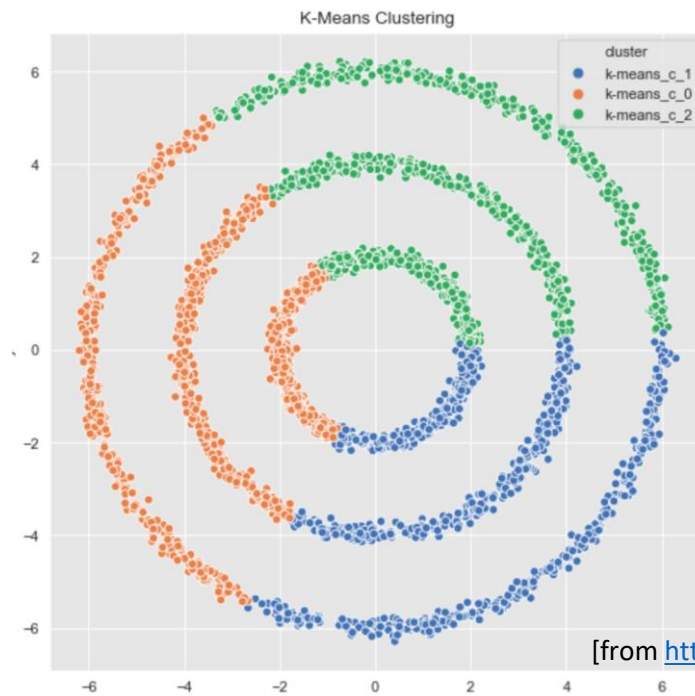
- The clusters returned by (hard) k-means algorithm are each convex.
- So k-means algo. won't work on a dataset where some "ground-truth" clusters are not convex.

- **Definition:**

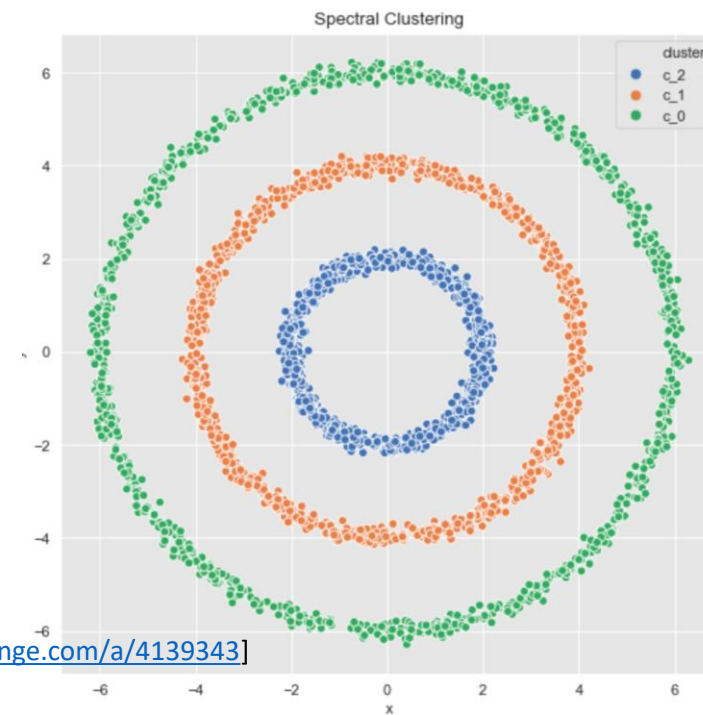
- Consider a dataset $D = \{x_1, x_2, \dots, x_n\}$ with each $x_i \in \mathbb{R}^d$, and a clustering of this dataset.
- A cluster C in this clustering is convex if for all pairs of points $x_i, x_j \in C$, any data point $x_k \in D$ that lie in the line connecting x_i, x_j also belongs to the same cluster C (i.e., $\forall x_k \in D$ that can be expressed as $\lambda x_i + (1 - \lambda)x_j$ for some $\lambda \in [0, 1]$, $x_k \in C$ also).

- **Examples:**

Convex clusters of D



Non-convex clusters (green and orange) of D



[from <https://math.stackexchange.com/a/4139343>]

Some final thoughts on Clustering and other
Clustering algorithms!

Clustering: a popular PR task

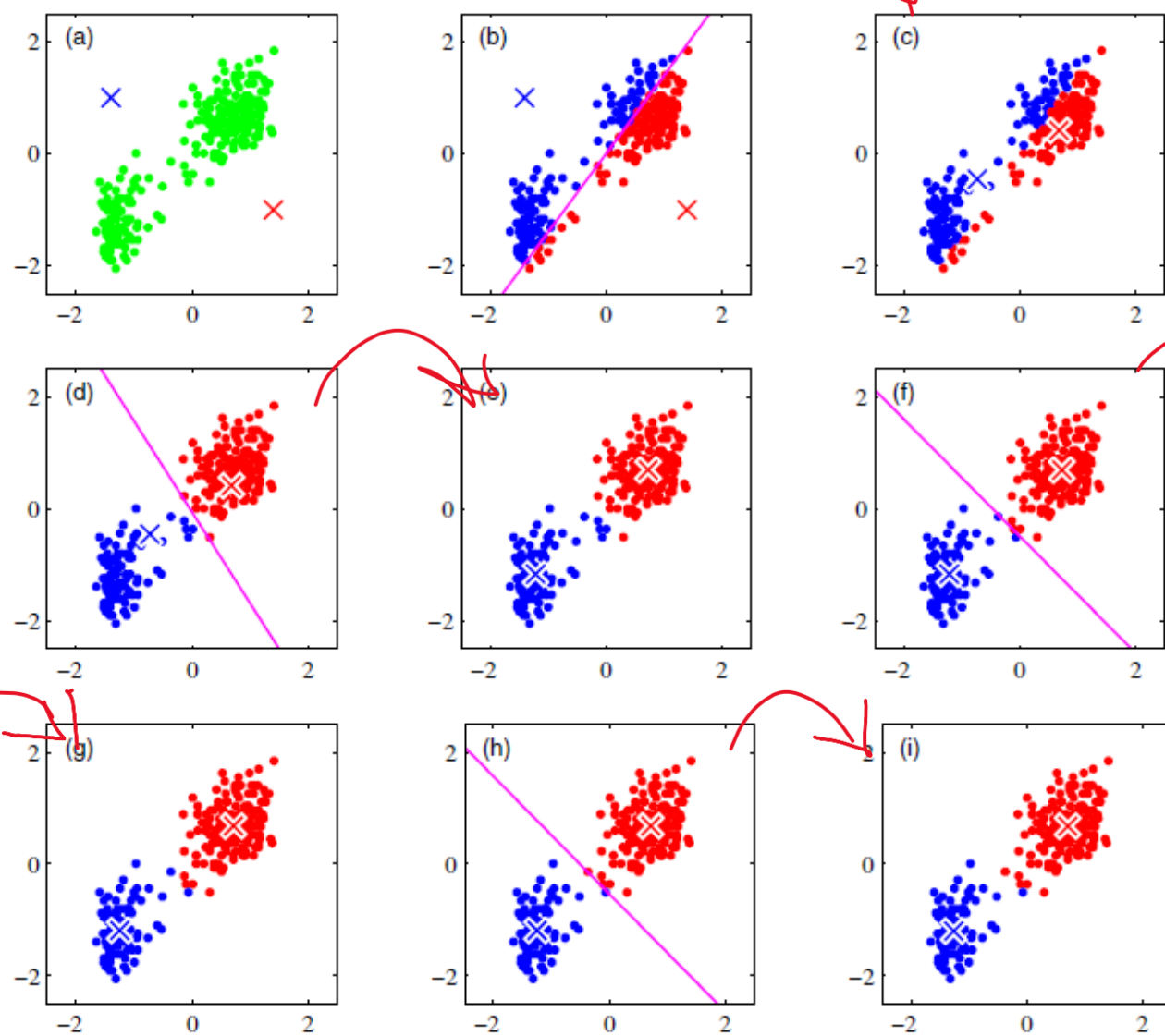
- Human brains are good at grouping objects based on their similarity. How do we automate it?
 - Group N objects into K groups
 - Paradigms: **K-means** clustering – **centroid-based** clustering; **Hierarchical** clustering – **nested/stratified clustering** with top-down or bottom-up approaches; **Spectral clustering** – **graph connectivity-based (in-syllabus, already seen)**, etc.
- Why is a good clustering popular? It
 - allows exploratory data analysis (unsupervised ML)
 - has predictive power
 - allows lossy compression
 - can reveal interesting outliers
- Our approach in these lectures
 - Focus on spectral clustering to illustrate applications of Linear Algebra in clustering and dim. redn. (as well as the connection between graph-connectivity and linear algebra)
 - Say a few words about k-means and hierarchical clustering, which you may've already seen in other courses, unlike spectral clustering

K-means clustering (in one slide! actually three :)

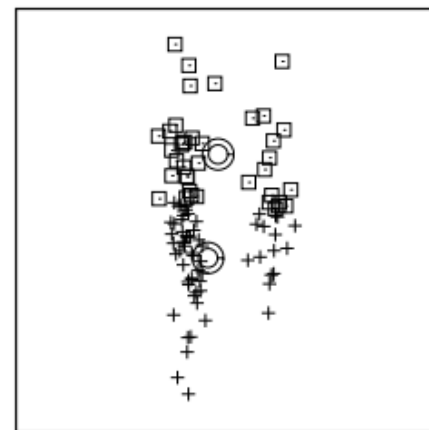
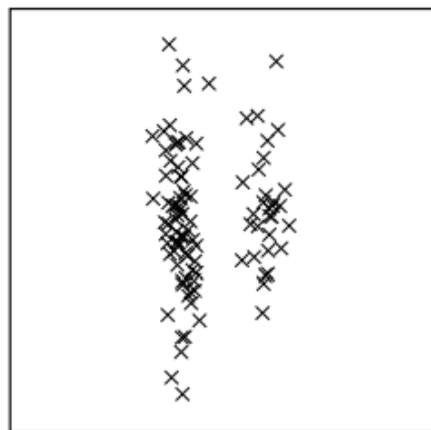
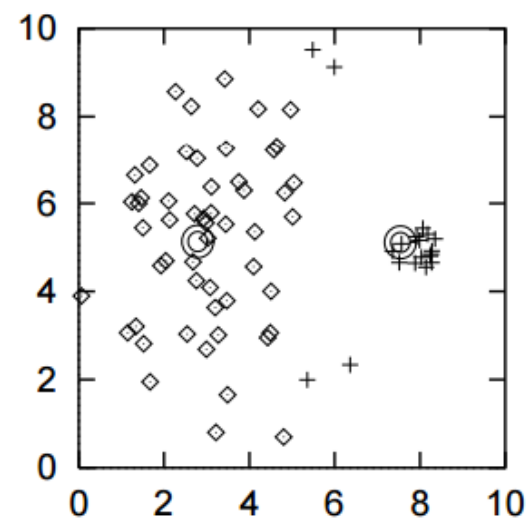
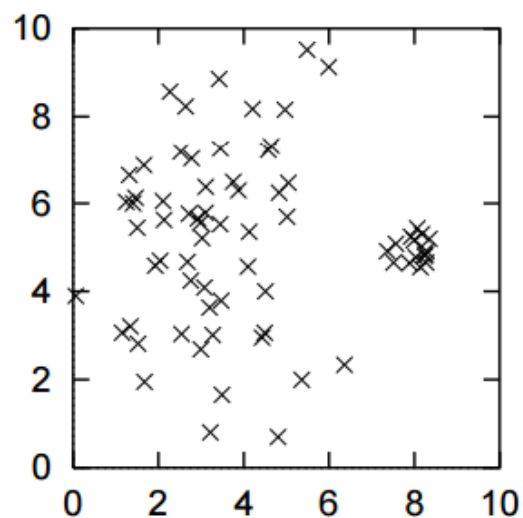
Problem: Find K cluster centers that minimize the sum of squared distance of each data point to the nearest cluster center.

Algorithm: Starting with K centers initialized in some way, iterate until convergence:

- a) [Centers -> Clusters] Assign* each data point to the nearest center.
- b) [Clusters -> Centers] Update* centers to sample means of data points they are responsible for.

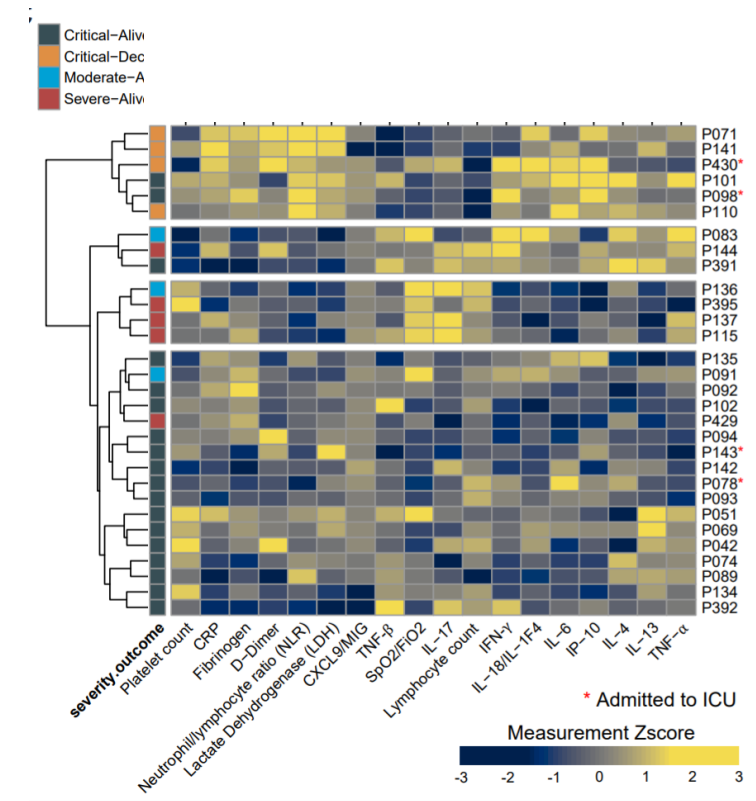


(hard) K-means pitfalls (and solutions?)



Hierarchical clustering (in one slide)

- Stratified clustering of objects, represented as a tree or dendrogram
 - very popular in heatmap display of data matrix.
 - E.g., in COVID bioinformatics shown to right!
- Two approaches to hierarchical clustering:
 - repeated bottom-up groupings of closest pair of objects – **agglomerative**
 - repeated top-down splitting of objects into groups – **divisive**
- In agglomerative, how do you construct the tree?
 - Convert the dataset into an adjacency or distance matrix as in Spectral clustering!
 - Aggregate the most similar pair of clusters/objects into one cluster, and iterate!



(Image: [Liu, Martins, Lau, et al. *Cell*. 2021])

Clustering: Overall Summary & Final Notes

- **Clustering** is a popular PR task that groups similar objects into clusters.
 - Unsupervised ML problem as it doesn't require any label information on the objects represented as datapoints.
 - Many algorithms available; the ones we saw were based on pairwise similarity/connectivity relations inferred from the datapoints: **K-means, Hierarchical, and Spectral**.
- In practice, evaluation of a clustering must be done via internal or external performance metrics:
 - **Internal:** how consistent is your clustering? E.g., Silhouette coefficient, Dunn Index
 - **External:** how well does your clustering match ground-truth labelling? E.g., Rand Index, F1-score.

pairwise L_2 \rightarrow +ive
 \rightarrow -ive

$$\frac{\min \text{ interclust dist}}{\max \text{ intracust dist}}$$