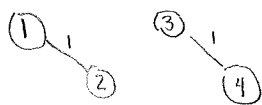


SPECTRAL CLUSTERING

Ex. TWO-ARMED SPIRAL AND CONCENTRIC RINGS

SEE CANVAS > FILES > NUMERICAL EXAMPLES > CLUSTERING

Ex: CONSIDER THE FOLLOWING SIMPLE GRAPH



$A =$ ADJACENCY (AFFINITY) MATRIX

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

A_{uv} = WEIGHT OF EDGE
CONNECTING NODES
 u, v

$D =$ DEGREE MATRIX

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

$D_{uu} = \sum_v A_{uv}$
= # NODES CONNECTED
TO u

$$L = \text{UNNORMALIZED GRAPH LAPLACIAN} = D - A = \begin{bmatrix} 1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 \\ 0 & 0 & 1 & -1 \\ 0 & 0 & -1 & 1 \end{bmatrix}$$

OBSERVE:

- L IS SYMMETRIC (DIAGONALIZABLE)

- GIVEN VECTOR \vec{x} , $\vec{x}^T L \vec{x} = \vec{x}^T D \vec{x} - \vec{x}^T A \vec{x}$
$$= \frac{1}{2} \sum_i \sum_j A_{ij} (x_i - x_j)^2 \geq 0 \text{ SINCE } A_{ij} \geq 0$$

$\therefore L$ IS POSITIVE-SEMIDEFINITE

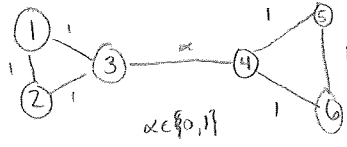
- 0 IS AN EIGENVALUE w/ EIGENVECTOR $\mathbf{1}_4$ (ALSO HAS EIGENVECTOR $(1, 1, 0, 0)^T$)

- 0 IS EIGENVALUE w/ GEOMETRIC/ALGEBRAIC MULTIPLICITY 2 = # CONNECTED COMPONENTS!

- UNCHANGED IF WE PERMUTE ROWS/COLUMNS (RE-LABEL THE NODES)

BUT KEEP THE SAME CONNECTIONS/STRUCTURE

Ex:



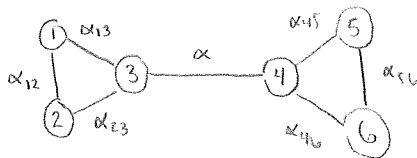
$$L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 & 0 \\ -1 & -1 & 2+\alpha & -\alpha & 0 & 0 \\ 0 & -\alpha & 2+\alpha & -1 & -1 & 0 \\ 0 & 0 & 0 & 2 & -1 & -1 \\ -1 & -1 & 0 & -1 & -1 & 2 \end{bmatrix}$$

- L STILL SYMMETRIC, POSITIVE-SEMI-DEFINITE
- 0 IS AN EIGENVALUE

EIGENVECTORS	$\alpha =$
$\mathbf{1}_4$	0
$(1, 1, 1, 0, 0, 0)^T$	1
$(0, 0, 0, 1, 1, 1)^T$	

- If $\alpha = 1$, $\lambda_2 = 0.4384$ IS THE SECOND SMALLEST EIGENVALUE
- λ_2 APPROXIMATES THE SPARSEST CUT OF THE GRAPH, i.e.
 THE RATIO OF THE FEWEST # OF REMOVED EDGES
 (WHICH SPLITS THE GRAPH INTO 2 DISCONNECTED PIECES)
 DIVIDED BY THE # OF VERTICES IN THE SMALLER OF
 THE 2 GRAPHS ($1/3$ IN EXAMPLE WHEN $\alpha = 1$)
- SPARSEST CUT RELATED TO CLUSTERING NODES INTO
 $K=2$ CLUSTERS WITH MINIMAL REMOVAL OF EDGES

Ex. WEIGHTED GRAPH



$$L = \begin{bmatrix} \alpha_{12} + \alpha_{13} & -\alpha_{12} & -\alpha_{13} & 0 & 0 & 0 \\ -\alpha_{12} & \alpha_{12} & 0 & 0 & 0 & 0 \\ -\alpha_{13} & 0 & \alpha_{13} + \alpha & -\alpha & 0 & 0 \\ 0 & 0 & -\alpha & \alpha + \alpha_{45} + \alpha_{46} & -\alpha_{45} & -\alpha_{46} \\ 0 & 0 & 0 & -\alpha_{45} & \alpha_{45} + \alpha_{56} & -\alpha_{56} \\ 0 & 0 & 0 & -\alpha_{46} & -\alpha_{56} & \alpha_{46} + \alpha_{56} \end{bmatrix}$$

- RELATIONSHIP BETWEEN λ_2 AND α UNCHANGED ($\lambda_2 = 0$ IFF $\alpha = 0$)

* SPECTRUM OF L IS GIVING INFORMATION ABOUT
 # OF CONNECTED COMPONENTS OF GRAPH!

ASIDE: WHY IS IT CALLED THE GRAPH LAPLACIAN?

SUPPOSE $\phi: V \longrightarrow \mathbb{R}^V$ DESCRIBES THE DISTRIBUTION OF HEAT
 \uparrow SET OF NODES \uparrow HEAT IN NODES

ACROSS A GRAPH $G = (V, E)$; $\phi_i(t)$ IS THE HEAT OF NODE i AT TIME t .

NEWTON'S LAW OF COOLING

THE HEAT TRANSFERRED BETWEEN NODES i AND j
 IS PROPORTIONAL (W/ CONSTANT $-A_{ij}$) TO $\phi_i - \phi_j$,
 THE DIFFERENCE IN HEAT BETWEEN NODES

$$\begin{aligned} \frac{d\phi_i(t)}{dt} &= - \sum_j A_{ij} (\phi_i(t) - \phi_j(t)) \\ &= - \phi_i \sum_j A_{ij} - \sum_j A_{ij} \phi_j \\ &= - \phi_i D_{ii} - \sum_j A_{ij} \phi_j = - \sum_j L_{ij} \phi_j \end{aligned}$$

IN MATRIX FORM

$$\frac{d}{dt} \phi = -L \phi$$

THIS IS THE SAME AS THE HEAT EQUATION W/ Δ , THE LAPLACIAN,
 REPLACED BY $-L$, HENCE "GRAPH LAPLACIAN"

THM: LET $G = (V, E)$ BE A WEIGHTED GRAPH, AND LET $L = D - A$
 BE ITS GRAPH LAPLACIAN. THEN

- i) THE NUMBER c OF CONNECTED COMPONENTS $K_1 \dots K_c$, IS EQUAL
 TO THE DIMENSION OF THE NULLSPACE OF L , WHICH IS
 EQUAL TO THE GEOMETRIC MULTIPLICITY OF $\lambda = 0$.
- ii) THE NULLSPACE HAS A BASIS CONSISTING OF INDICATOR VECTORS
 OF THE CONNECTED COMPONENTS OF G , THAT IS VECTORS $(v_1 \dots v_N)^T$
 S.T. $v_j = 1$ IFF NODE j IS IN COMPONENT K_i , 0 OTHERWISE. FOR $i = 1 \dots c$

DATA SIMILARITY : GRAPH LAPLACIAN

Given data $\vec{x}_1, \dots, \vec{x}_n$ construct an $N \times N$ affinity matrix, A ,

so that $A_{ij} = A_{ji}$ is the similarity between \vec{x}_i and \vec{x}_j .

* Common choice is $A_{ij} = \exp\left(-\frac{\|\vec{x}_i - \vec{x}_j\|_2}{2\sigma^2}\right)$
GAUSSIAN RADIAL
BASIS FUNCTION

EQUIVALENTLY

$$\exp(-\gamma \|\vec{x}_i - \vec{x}_j\|_2)$$

$$\exp(-\|\vec{x}_i - \vec{x}_j\|_2^2 / c)$$

Then create GRAPH LAPLACIAN

$$L = D - A \quad , \quad \text{NOTE } A_{ii} = 1 \text{ BUT CANCELS IN } L$$

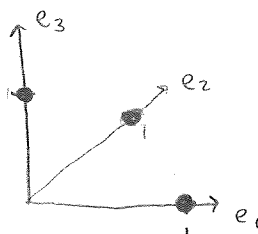
$$D_{ii} = \sum_j A_{ij} = 1 + \sum_{j \neq i} A_{ij}$$

$$L_{ii} = D_{ii} - A_{ii} = \sum_{j \neq i} A_{ij}$$

IDEAL CASE:

- $\|\vec{x}_i - \vec{x}_j\| = \infty$ WHEN \vec{x}_i, \vec{x}_j IN DIFFERENT COMPONENTS
- L HAS EIGENVALUE 0 W/ MULTIPLICITY K
- $\vec{y}_1, \dots, \vec{y}_k \in \mathbb{R}^N$ LINEARLY INDEPENDENT EIGENVECTORS w/ EIGENVALUE 0
- TREAT $Y = \begin{bmatrix} \vec{y}_1 & \dots & \vec{y}_k \end{bmatrix} \in \mathbb{R}^{N \times k}$ AS A DATA MATRIX

WHAT DOES GRAPH OF (ROWS) OF Y LOOK LIKE?



ALL POINTS IN CLUSTER i
MAPPED TO \vec{e}_i ! (IDEALLY?)

APPLYING K-MEANS CLUSTERING HERE WORKS VERY, VERY WELL !!

NOTE: EVEN IN THE IDEAL CASE IT IS UNLIKELY THAT EIGENVALUE SOLVER WILL GIVE

$$Y = \begin{bmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 0 \end{bmatrix}$$

HOWEVER, $\vec{y}_1, \dots, \vec{y}_k$ WILL BE ORTHOGONAL (ORTHONORMAL IF NORMALIZED)

BY THE SPECTRAL DECOMPOSITION THM. (LIMITS THE INTERPRETABILITY OF Y)

\Rightarrow CAN CREATE K DISTINCT CLUSTERS ON THE UNIT SPHERE IN \mathbb{R}^k WHICH ARE WELL SEPARATED!

\Rightarrow APPLYING K -MEANS CLUSTERING WILL THEN CAPTURE THE K CLUSTERS IN OUR DATA!

* IN CASES OF EXTREME σ OR $\|x_i - x_j\|$, TRADITIONAL EIGENVALUE SOLVERS CAN STRUGGLE BUT L IS POSITIVE SEMIDEFINITE SO THAT USING SVD GIVES MORE RELIABLE RESULTS

THE UNNORMALIZED CASE (AKA REALITY)

• MANY, MANY CHOICES FOR CONSTRUCTING THE AFFINITY MATRIX A

• $\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$ IS JUST ONE OF MANY DIFFERENT "KERNELS"

THAT ONE CAN USE. MORE ON THIS NEXT WEEK.

• STILL NEED TO TUNE σ^2 IN RADIAL BASIS FUNCTION, $\exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$

• CAN ALSO CHOOSE TO ONLY USE NEIGHBORHOOD INFORMATION (HASTIE TEXT)

i.e. SET $A_{ij} = 0$ IF i, j ARE NOT K -NEAREST NEIGHBORS

WHICH GIVES A SPARSE AFFINITY MATRIX

• UNNORMALIZED: $L = D - A$

SYMMETRIC NORMALIZED: $L^{sym} = D^{-1/2} L D^{-1/2} = I - D^{-1/2} A D^{-1/2}$

RANDOM WALK NORMALIZED: $L^{rw} = D^{-1} L = I - D^{-1} A$

• ALSO

• # OF EIGENVALUES OF L TO USE (NORMALIZED)

• # OF CLUSTERS

CASE STUDY: ON SPECTRAL CLUSTERING (AL, JORDAN, WEISS)

- * AUTOMATIC FORM OF $L = D^{1/2} A D^{1/2}$ SEPARATES $\lambda \mapsto 1/\lambda$ FROM OUR CONNECTION
- * IMPLEMENTED IN SPECT FUNCTION

ALGORITHM

GIVEN $\vec{x}_1, \dots, \vec{x}_n \in \mathbb{R}^d$ THAT WE WANT TO CLUSTER INTO k SUBSETS

1) FORM AFFINITY $A \in \mathbb{R}^{n \times n}$

$$A_{ij} = \exp\left(-\|\vec{x}_i - \vec{x}_j\|^2 / 2\sigma^2\right) \quad i \neq j, \quad A_{ii} = 0$$

2) DEFINE $D \in \mathbb{R}^{n \times n}$ TO BE DIAGONAL MATRIX w/ $D_{ii} = \sum_j A_{ij}$

also let $L = D^{-1/2} A D^{-1/2}$

3) FIND THE k -LARGEST EIGENVALUES AND THEIR ASSOCIATED

EIGENVECTORS (CHOSEN TO BE ORTHOGONAL IF EIGENVALUES ARE

REPEATED) $\vec{w}_1, \dots, \vec{w}_k \in \mathbb{R}^n$ AND FORM $W = [\vec{w}_1, \dots, \vec{w}_k]$

4) FORM $Y \in \mathbb{R}^{n \times k}$ FROM W BY

$$Y_{ij} = \frac{W_{ij}}{\left(\sum_j W_{ij}^2\right)^{1/2}} \quad \left(\text{PUT ROWS OF UNIT SPHERE}\right)$$

5) TREAT ROWS OF Y AS POINTS IN \mathbb{R}^k , CLUSTER INTO k -CLUSTERS VIA k -MEANS

6) ASSIGN \vec{x}_i TO CLUSTER j IFF ROW i OF Y WAS ASSIGNED TO CLUSTER j .

* ANALYSIS OF AL, JORDAN, AND WEISS GIVES A LESS AD HOC METHOD FOR CHOOSING σ !

* SEARCH OVER A RANGE OF σ^2 (OBTAINED BY LOOKING AT VARIABILITY IN A SUBSET OF DATA) AND CHOOSE VALUE GIVING THE TIGHTEST CLUSTERS! (LOWEST WITHIN CLUSTER SUM OF SQUARES)

* CAN COMPARE PERFORMANCE BY LOOKING AT TOTAL OF WITHINSS AS GLOBAL MEASURE OF ERROR

Ex: SPIRAL & CONCENTRIC RINGS

CASUAL (11-CLASS EXERCISE)