

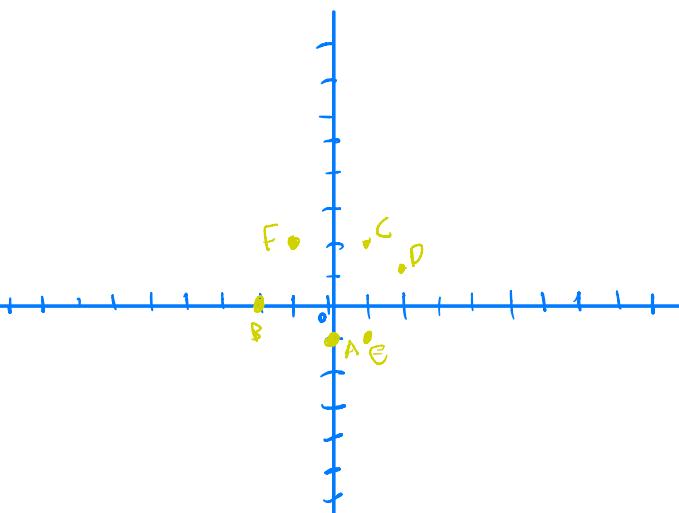
CZ/CE 4032
Data Analytics & Mining

Tutorial – Week 2
Clustering – basic methods

Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

Questions

- Q1:
- Perform K-means on the following 6 data points for 2 iterations. The number of clusters is 2. The initial centroids are at data point E and F.



A	[0, -1]
B	[-2, 0]
C	[1, 2]
D	[2, 1]
E	[1, -1]
F	[-1, 2]

Questions

- Q2:
- Given the 6 data points in Q1, perform K-means for 2 iterations. The number of clusters is 2. The initial centroids are at data point C and A.

	$C(1, 2)$	$A(0, -1)$
B(-2, 0)	$q+4 = 13$	$4+1 = 5$
D(2, 1)	$1+1 = 2$	$4+4 = 8$
E(1, -1)	9	1
F(-1, 2)	4	$1+9 = 10$

$$C_1 = (0, F)$$
$$C = \left(\frac{1}{2}, \frac{3}{2} \right)$$

$$C_2 = (B, E)$$
$$A = \left(-\frac{1}{2}, -\frac{1}{2} \right)$$

Iteration 2

$$\begin{cases} C(1/2, 3/2) \\ \begin{aligned} B(-2, 0) & \frac{28}{4} + \frac{9}{4} = \frac{37}{4} \\ D(2, 1) & \frac{9}{4} + \frac{4}{4} = \frac{13}{4} \\ E(1, -1) & \\ F(-1, 2) & \end{aligned} \end{cases}$$

$$A = (-1/2, -1/2)$$

Questions

- Q3:
- Given the 6 data points in Q1, perform hierarchical clustering (agglomerative) until there is only 1 cluster. Use Single Linkage method.

Iteration 1
step 2 = Cluster assignment
3) Calculate squared L2 distance

	$E(1, -1)$	$F(-1, 2)$
A	$(0, -1)$	$(0+1)^2 + (-1+1)^2 = 1$
B	$(-2, 0)$	$(-2+1)^2 + (0+1)^2 = 5$
C	$(1, 2)$	$(1-1)^2 + (2+1)^2 = 9$
D	$(2, 1)$	$(2-1)^2 + (1-1)^2 = 1$
E	$(1, -1)$	$4+9 = 13$
F	$(-1, 2)$	$4+9 = 13$

2) Assign each point to its nearest centroid

$$C_1 : (A, D, E) \quad C_2 : (B, C, F)$$

$$C_3 : \left(\frac{0+2+1}{3}, \frac{-1+1-1}{3} \right) = \left(\frac{3}{3}, \frac{0+2+2}{3} \right) = \left(1, \frac{4}{3} \right)$$

Iteration 2

	$(1, -\frac{1}{3})$	$(-\frac{2}{3}, \frac{4}{3})$
A	$(0, -1)$	$1+\frac{1}{9} = \frac{10}{9} = 1.11$
B	$(-2, 0)$	$9+\frac{4}{9} = \frac{85}{9} = 9.44$
C	$(1, 2)$	$16/9 + 16/9 = \frac{32}{9} = 3.56$
D	$(2, 1)$	$25/9 + 4/9 = \frac{29}{9} = 3.22$
E	$(1, -1)$	$4/9 + 1/9 = \frac{5}{9} = 0.56$
F	$(-1, 2)$	8.22

$$2) \quad C_1(C, D, F) \quad C_2(A, B, E)$$

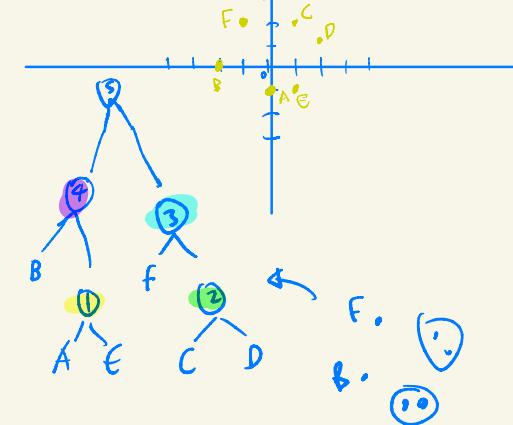
A	$(0, -1)$	10	0
B	$(-2, 0)$	13	5
C	$(1, 2)$	0	10
D	$(2, 1)$	2	8
E	$(1, -1)$	9	1
F	$(-1, 2)$	4	10

$$C_1(C, D, F) = \left(\frac{1+2-1}{3}, \frac{2+1+2}{3} \right) = \left(\frac{2}{3}, \frac{5}{3} \right)$$

$$C_2(A, B, E) = \left(\frac{0-2+1}{3}, \frac{-1+0+1}{3} \right) = \left(-\frac{1}{3}, -\frac{1}{3} \right)$$

A	$(0, -1)$	7.56	0.22
B	$(-2, 0)$	9.89	3.22
C	$(1, 2)$	0.22	8.89
D	$(2, 1)$	2.22	0.22
E	$(1, -1)$	7.22	1.89
F	$(-1, 2)$	2.89	7.56

3)



$$\begin{array}{lllll} A(0, -1) & B(-2, 0) & C(1, 2) & D(2, 1) & E(1, -1) & F(-1, 2) \\ 0 & 4+1=5 & 1+9=10 & 4+4=8 & 1+0=1 & 1+9=10 \\ 0 & 9+9=18 & 16+1=17 & 9+1=10 & 1+4=5 & 1+4=5 \\ 0 & 1+1=2 & 0+9=9 & 4+0=4 & 0 & 4+0=4 \\ 0 & 1+9=10 & 9+1=10 & 0 & 0 & 7+1=10 \\ 0 & & & 0 & & 7+9=16 \end{array}$$

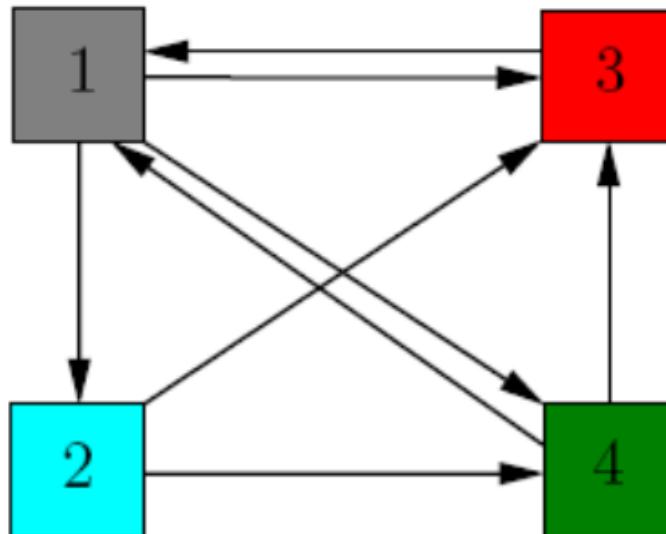
$$\begin{array}{l} B(-2, 0) \\ AE \\ CDF \end{array}$$

! we have 3 equal distances,
we can randomly select

Link Analysis: PageRank Tutorial

Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

- Q1. Given the graph below, calculate the page rank score for each node using the following two methods:
 - 1) directly solving the flow equations;
 - 2) using the power iterative method for 2 iterations.



1) Solve

$$\begin{array}{c|cc} & \text{in} & \text{out} \\ \hline 1 & 2 & 3 \\ 2 & 1 & 2 \\ 3 & 3 & 1 \\ 4 & 2 & 2 \end{array}$$

$$r_1 + r_2 + r_3 + r_4 = 1 \quad (1)$$

$$r_1 = r_3 + \frac{1}{2}r_4 \quad (2)$$

$$r_2 = \frac{1}{3}r_3 \quad (3)$$

$$r_3 = \frac{1}{3}r_1 + \frac{1}{2}r_2 + \frac{1}{2}r_4 \quad (4)$$

$$r_4 = \frac{1}{3}r_1 + \frac{1}{3}r_2 \quad (5)$$

use (2) in (4)

$$r_4 = \frac{1}{3}r_1 + \frac{1}{2}\left(\frac{1}{3}r_1\right) = \frac{1}{2}r_1 \quad (6)$$

use (6) in (1)

$$r_1 = r_3 + \frac{1}{2}\left(\frac{1}{2}r_1\right) = r_1 - \frac{1}{4}r_1 = r_3$$

$$r_3 = \frac{3}{4}r_1 \quad (7)$$

Transition matrix

$$\begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 2 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 3 & 1 & 0 & 0 & 0 \\ 4 & \frac{1}{3} & 0 & \frac{1}{2} & 0 \end{pmatrix}$$

$\xrightarrow{\text{in out}}$ by T o T is $\xrightarrow{\text{out in}}$

use (6) (5) (2) in (1)

$$r_1 + \left(\frac{1}{3}r_1\right) + \left(\frac{3}{4}r_1\right) + \left(\frac{1}{2}r_1\right) = 1$$

$$\frac{12 + 4 + 9 + 6}{12} r_1 = 1$$

$$r_1 = \frac{12}{31}$$

$$r_2 = \frac{4}{31}$$

$$r_3 = \frac{9}{31}$$

$$r_4 = \frac{6}{31}$$

2) $\Gamma = M \cdot \Gamma'$

Power iteration

Transition matrix

$$\begin{pmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{pmatrix}$$

$$\Gamma_0 = \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix}$$

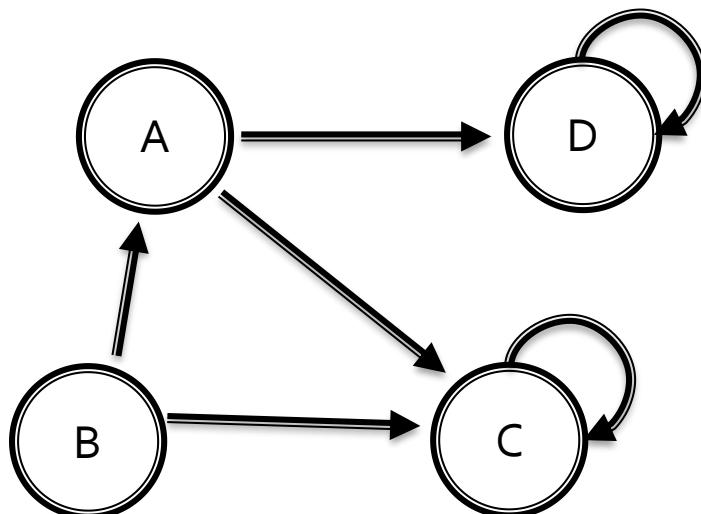
$$\Gamma_1 = M \cdot \Gamma$$

$$\Gamma_1 = \begin{pmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix} = \begin{pmatrix} \frac{1}{4} + \frac{1}{8} \\ \frac{1}{12} \\ \frac{1}{12} + \frac{1}{8} + \frac{1}{8} \\ \frac{1}{12} + \frac{1}{8} \end{pmatrix} = \begin{pmatrix} \frac{3}{8} \\ \frac{1}{12} \\ \frac{1}{3} \\ \frac{5}{24} \end{pmatrix}$$

$$\Gamma_2 = M \cdot \Gamma_1$$

$$\Gamma_2 = \begin{pmatrix} 0 & 0 & 1 & \frac{1}{2} \\ \frac{1}{3} & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{3} & \frac{1}{2} & 0 & 0 \end{pmatrix} \cdot \begin{pmatrix} \frac{3}{8} \\ \frac{1}{12} \\ \frac{1}{3} \\ \frac{5}{24} \end{pmatrix} = \begin{pmatrix} \frac{1}{8} + \frac{3}{48} \\ \frac{3}{12} \\ \frac{3}{12} + \frac{1}{24} + \frac{5}{48} \\ \frac{3}{12} + \frac{1}{24} \end{pmatrix} = \begin{pmatrix} \frac{7}{48} \\ \frac{1}{8} \\ \frac{17}{48} \\ \frac{1}{6} \end{pmatrix} = \begin{pmatrix} 0.1458 \\ 0.125 \\ 0.3542 \\ 0.1667 \end{pmatrix}$$

- Q2. A graph is given below.
 - a) calculate the page rank score for each node using the power iteration method for 3 iterations.
 - b) identify one spider trap group if there are any.



	in	out
A	1	2
B	0	2
C	3	1
D	2	1

Matrix =

$$\Gamma_0 = \begin{pmatrix} \frac{1}{4} & & & \\ & \frac{1}{4} & & \\ & \frac{1}{4} & & \\ & \frac{1}{4} & & \end{pmatrix}$$

$$\Gamma_0 = \left(\begin{array}{cccc} 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 0 & 0 & 1 \end{array} \right)$$

$$\Gamma_{t+1} = M \cdot \Gamma_t$$

$$\Gamma_1 = \begin{pmatrix} 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \\ \frac{1}{4} \end{pmatrix} = \begin{pmatrix} \frac{1}{8} \\ 0 \\ \frac{1}{8} + \frac{1}{8} + \frac{1}{4} \\ \frac{1}{8} + \frac{1}{4} \end{pmatrix} = \begin{pmatrix} \frac{1}{8} \\ 0 \\ \frac{3}{8} \\ \frac{3}{8} \end{pmatrix}$$

$$\Gamma_2 = \begin{pmatrix} 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{8} \\ 0 \\ \frac{3}{8} \\ \frac{3}{8} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \frac{1}{16} + \frac{1}{2} \\ \frac{1}{16} + \frac{3}{8} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \frac{9}{16} \\ \frac{7}{16} \end{pmatrix}$$

$$\Gamma_3 = \begin{pmatrix} 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ \frac{9}{16} \\ \frac{7}{16} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \frac{9}{16} \\ \frac{7}{16} \end{pmatrix}$$

Only Spider trap node have Values

3) Spider trap = no out link node
 Spider trap groups =
 = C - A, C, D = Ac?
 = D - C, D = Bc?

Spider trap is not dead node

!! deadnode is totally no out-link
 Spider trap might exist with node only linking to itself

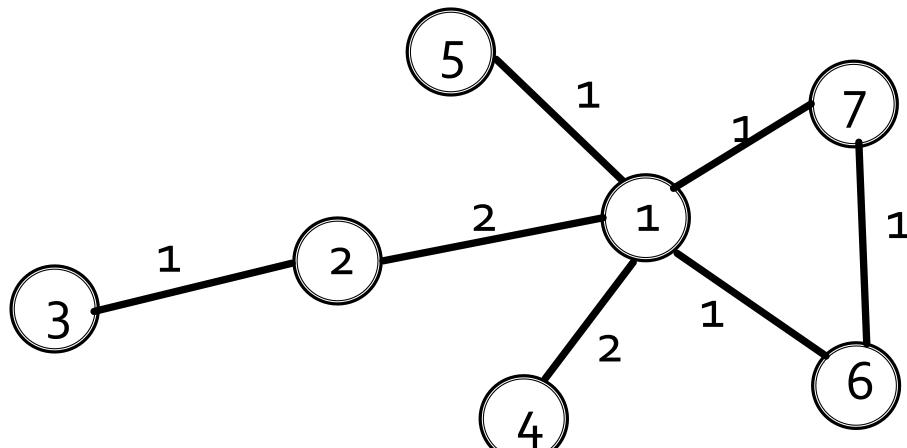


Graph Neural Network Tutorial

Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

Question:

Given a graph below (next page), the task is to do node-wise classification using a 2-layer graph convolutional network (GCN) and a cross-entropy loss.



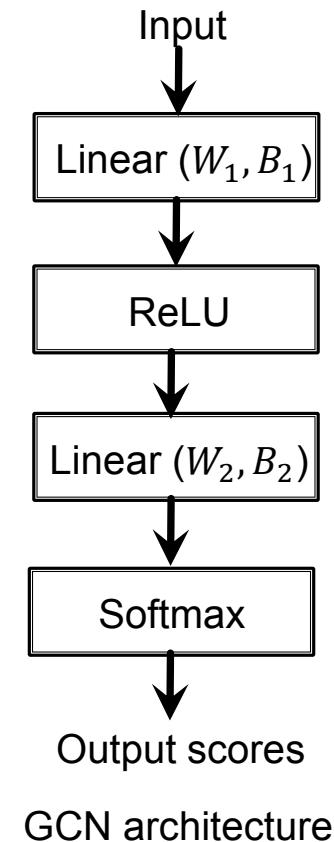
Input graph

need calculate node 3 and 2

(continued from last page)

The initial node features and the GCN network architecture are given below:

$$\begin{aligned}x_1 &= [0, -1]^T, \\x_2 &= [0, 1]^T, \\x_3 &= [-1, 0]^T, \\x_4 &= [1, 0]^T, \\x_5 &= [-1, -1]^T, \\x_6 &= [2, 1]^T, \\x_7 &= [-1, 1]^T.\end{aligned}$$



(continued from last page)

The initial GCN weight parameters given below (W_k and B_k are the weight matrices for neighborhood aggregation and self transformation, respectively, for the k -th layer).

Calculate the prediction of **node 3** by performing one forward pass.

$$W_1 = \begin{bmatrix} 0 & -0.1 \\ 0.1 & 0 \end{bmatrix} \quad B_1 = \begin{bmatrix} 0.1 & 0 \\ 0 & -0.1 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} -0.1 & 0 \\ 0 & 0.1 \end{bmatrix} \quad B_2 = \begin{bmatrix} 0 & 0.1 \\ -0.1 & 0 \end{bmatrix}$$

1) input
 \downarrow
 linear
 \downarrow
 ReLU
 \downarrow
 Linear
 \downarrow
 Softmax

$$\hat{h}_i^{(k)} = w_k \sum_{j \in N_i} a'_{ij} h_j^{(k-1)} + b_k h_i^{(k-1)}$$

$$\rightarrow \max(0, h)$$

row-normalized
 affinity matrix

$$\begin{matrix} 0 & 2 & 0 & 2 & 1 & 1 & 1 \\ 2 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 2 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \end{matrix}$$

$$\begin{matrix} 0 & 2/4 & 0 & 2/4 & 1/4 & 1/4 & 1/4 \\ 2/3 & 0 & 1/3 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/3 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 0 & 0 & 0 & 0 & 1/2 \\ 1/2 & 0 & 0 & 0 & 0 & 1/2 & 0 \end{matrix}$$

node 3 by fw pass

$$\begin{aligned} \text{node 1} \\ \text{Let } a'_{12} h_2 + a'_{14} h_4 + a'_{15} h_5 + a'_{16} h_6 + a'_{17} h_7 \\ \begin{pmatrix} 0 \\ 1 \\ -1 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} + \begin{pmatrix} -1 \\ 0 \\ -1 \end{pmatrix} + \begin{pmatrix} 2 \\ 1 \\ 1 \end{pmatrix} + \begin{pmatrix} -1 \\ 1 \\ -1 \end{pmatrix} = \end{aligned}$$

$$w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} =$$

$$\begin{aligned} \text{node 2} \\ \text{Let } a'_{21} h_1 + a'_{23} h_3 \\ \begin{pmatrix} 0 \\ -1 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \end{aligned}$$

$$w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} =$$

$$\begin{aligned} \text{node 3} \\ \text{Let } a'_{32} h_2 \\ \begin{pmatrix} 0 \\ 0 \end{pmatrix} = 0 \\ w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0 \end{pmatrix} \\ B_k h_3 = \begin{pmatrix} 0.1 & 0 \\ 0 & -0.1 \end{pmatrix} \begin{pmatrix} -1 \\ 0 \end{pmatrix} = \begin{pmatrix} -0.1 \\ 0 \end{pmatrix} \end{aligned}$$

$$\begin{aligned} \text{node 4} \\ \text{Let } a'_{41} h_1 \\ \begin{pmatrix} 0 \\ -1 \end{pmatrix} = \end{aligned}$$

$$w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} =$$

$$\begin{aligned} \text{node 5} \\ \text{Let } a'_{51} h_1 \\ w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} = \end{aligned}$$

$$\begin{aligned} \text{node 6} \\ \text{Let } a'_{61} h_1 + a'_{67} h_7 \\ \begin{pmatrix} 0 \\ -1 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \\ w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} = \end{aligned}$$

$$\begin{aligned} \text{node 7} \\ \text{Let } a'_{71} h_1 + a'_{76} h_6 \\ \begin{pmatrix} 0 \\ 0 \end{pmatrix} + \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \\ w_k \sum a'_{ij} h_j^{(k-1)} = \begin{pmatrix} 0 & -0.1 \\ 0.1 & 0 \end{pmatrix} = \end{aligned}$$

$$\text{Linear} = \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix} + \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix} = \begin{pmatrix} -0.2 \\ 0 \end{pmatrix}$$

$$\text{ReLu} = \max\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.1 \\ 0.1 \end{pmatrix}\right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \leftarrow h_2^{(1)}$$

↑
might be diff

Step 2 calculate $h_2^{(1)}$
 $| \quad \text{ReLU} = \begin{pmatrix} 0.067 \\ 0 \end{pmatrix} \leftarrow h_2^{(1)}$

GCN layer 2 !!

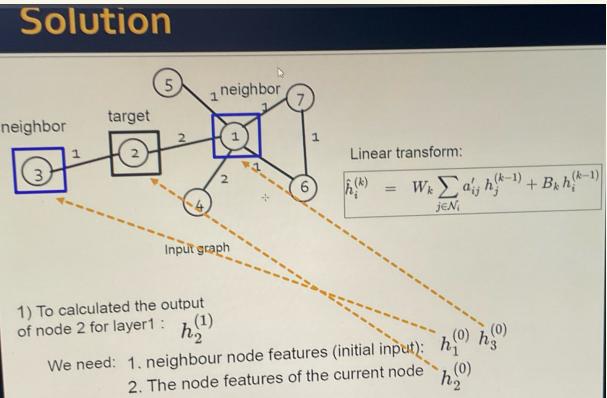
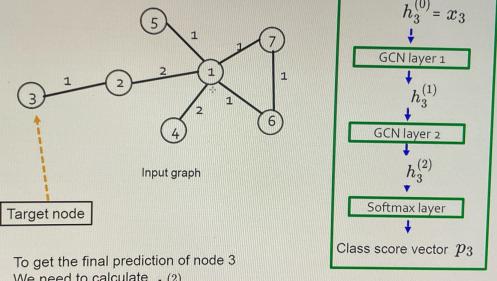
$h_3^{(2)}$

$$W_2 \in \alpha_{32} h_2 = \begin{pmatrix} -0.1 & 0 \\ 0 & 0.1 \end{pmatrix} \begin{pmatrix} 1. & 0.067 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} -0.0067 \\ 0 \end{pmatrix}$$

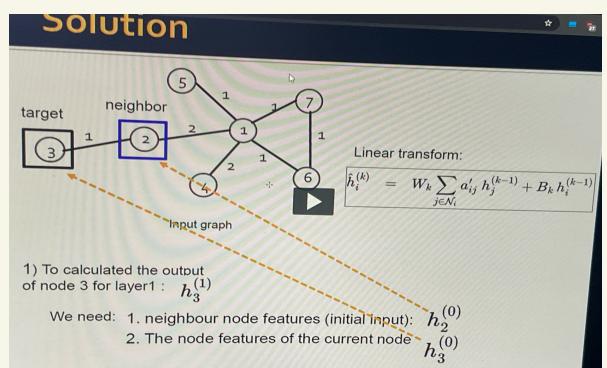
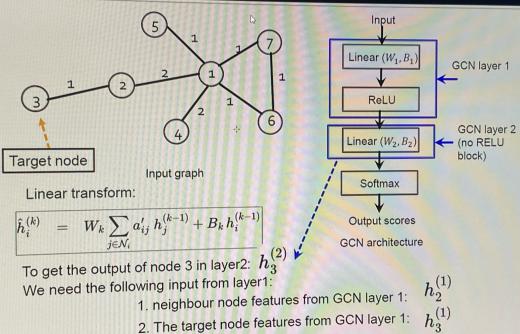
$$b_2 h_3 = \begin{pmatrix} 0 & 0.1 \\ -0.1 & 0 \end{pmatrix} \cdot \begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$h_3^{(2)} = \begin{pmatrix} -0.0067 \\ 0 \end{pmatrix} \rightarrow \text{Soft max} \rightarrow \begin{pmatrix} 0.418 \\ 0.582 \end{pmatrix}$$

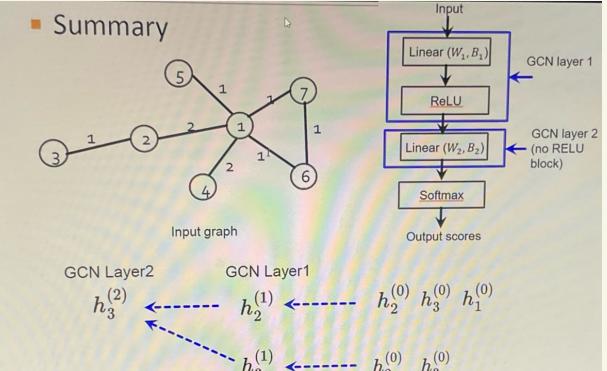
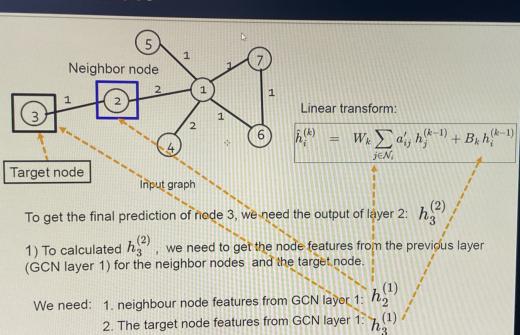
Solution



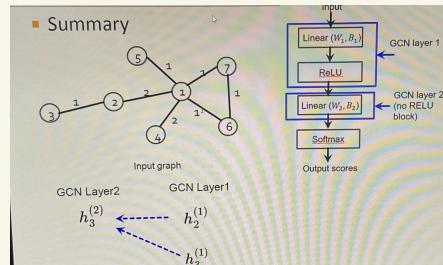
Solution



Solution



Summary



Similarity Search Tutorial

Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

Q1 LSH

- **Q1:** given the query sample and the samples in the database below, find the top 2 nearest neighbours in the database. Use LSH-random projection with the hash functions given below.

Query sample:

A	[0, -1]
---	----------

$$w_1 = [-1 \ 1]^T;$$

$$w_2 = [-1 \ 0]^T;$$

$$w_3 = [0 \ 1]^T;$$

$$w_4 = [1 \ -1]^T;$$

$$w_5 = [1 \ 0]^T;$$

$$w_6 = [-1 \ -1]^T;$$

Database samples (5):

B	[-2, 0]
C	[1, 2]
D	[2, 1]
E	[1, -1]
F	[-1, 2]

Note: the question is modified from the LSH example in the lecture class by adding two additional hashing functions (w_5 and w_6)

Query sample:

A	[0, -1]
---	-----------

$$\begin{aligned} w_1 &= [-1 \ 1]^T; \\ w_2 &= [-1 \ 0]^T; \\ w_3 &= [0 \ 1]^T; \\ w_4 &= [1 \ -1]^T; \\ w_5 &= [1 \ 0]^T; \\ w_6 &= [-1 \ -1]^T; \end{aligned}$$

Database samples (5):

B	[-2, 0]
C	[1, 2]
D	[2, 1]
E	[1, -1]
F	[-1, 2]

A $h_1 = w_1^T x_A = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = -1 \quad h_1(x_A) = 0$
 $h_2 = w_2^T x_A = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = 0 \quad h_2(x_A) = 0$
 $h_3 = w_3^T x_A = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = -1 \quad h_3(x_A) = 0$
 $h_4 = w_4^T x_A = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = 1 \quad h_4(x_A) = 1$
 $h_5 = w_5^T x_A = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = 0 \quad h_5(x_A) = 0$
 $h_6 = w_6^T x_A = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} 0 \\ -1 \end{pmatrix} = 1 \quad h_6(x_A) = 1$

B $h_1 = w_1^T x_B = 2 \quad h_1(x_B) = 1$
 $h_2 = \dots = 2 \quad h_2(x_B) = 1$
 $h_3 = \dots = 0 \quad h_3(x_B) = 0$
 $h_4 = \dots = 0 \quad h_4(x_B) = 0$
 $h_5 = \dots = 0 \quad h_5(x_B) = 0$
 $h_6 = \dots = 0 \quad h_6(x_B) = 0$

C $h_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = -2 \rightarrow 0$
 $h_2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = -1 \rightarrow 0$
 $h_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = -1 \rightarrow 0$
 $h_4 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 2 \rightarrow 1$
 $h_5 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 1$
 $h_6 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = -2 \rightarrow 0$

A	0	0	0	1	0
B	1	1	0	0	1
C	1	0	1	0	1
D	0	0	1	1	0
E	0	0	0	1	1
F	1	1	0	0	0

$$\begin{aligned} h_1 &= \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = -1+2 = 1 \\ h_2 &= \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = -1 \rightarrow 0 \\ h_3 &= \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 2 \rightarrow 1 \\ h_4 &= \begin{pmatrix} 1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 1-2 = -1 \rightarrow 0 \\ h_5 &= \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = 1+0 = 1 \\ h_6 &= \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = -3 \rightarrow 0 \end{aligned}$$

D $h_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = -1 \rightarrow 0$
 $h_2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = -2 \rightarrow 0$
 $h_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = 1$
 $h_4 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = 1$
 $h_5 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = 2 \rightarrow 1$
 $h_6 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = -3 \rightarrow 0$

E $h_1 = \begin{pmatrix} -1 \\ 1 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = 1+2 = 3 \rightarrow 1$
 $h_2 = \begin{pmatrix} -1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = 1$
 $h_3 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = 2 \rightarrow 0$
 $h_4 = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = -1-2 = -3 \rightarrow 0$
 $h_5 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = -1 \rightarrow 0$
 $h_6 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \begin{pmatrix} -1 \\ 2 \end{pmatrix} = 1-2 = 0$

A B C D E F
 A 3 S 3 2 4

L₂ distances

$$\begin{array}{llllll} A(0,-1) & B(-2,0) & C(1,2) & D(2,-1) & E(1,-1) & F(-1,2) \\ & 4+1=\underline{5} & 1+9=\underline{10} & 4+4=\underline{8} & 1+0=\underline{1} & 1+9=\underline{10} \end{array}$$

Q2 PQ

■ Q2. Product Quantization (PQ)

! Study asymmetry!

Two input vectors are given below:

$$A: [1, 0, 2, 1, -1, 0]$$

$$B: [-1, 3, 1, 4, 3, 2]$$

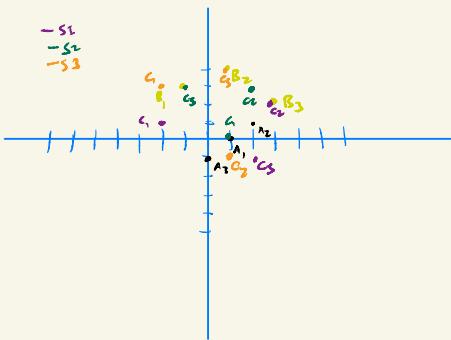
Define 3 Subspaces. Each subspace has 3 centroids:

$$\text{Subspace 1: } C1: [-2, 1], \ C2: [3, 2], \ C3: [2, -1]$$

$$\text{Subspace 2: } C1: [1, 0], \ C2: [2, 3], \ C3: [-1, 3]$$

$$\text{Subspace 3: } C1: [-2, 3], \ C2: [1, -1], \ C3: [1, 4]$$

- a): Use PQ to compute the compressed vectors of A and B
- b): Construct distance lookup tables and calculate the approximate (symmetric case) squared-L2 distance of A and B using the compressed vectors.



Squared L2 distance			
$A(-2, 1)$	$A(2, 1)$	$A(1, 0)$	$B(1, 2)$ $B(1, 1)$ $B(3, 2)$
$9+1=10$	$9+1=10$	$1+1=2$	$1+4+5=10$ $1+1=2$
$C(3, 2)$	$C(1, 2)$	$C(2, 1)$	$16+1=17$ $9+16=25$
$C(2, 1)$	$C(1, 0)$	$C(1, 1)$	$0+16=16$ $1+1=2$
$C(1, 0)$	$C(2, 3)$	$C(1, 1)$	$4+1=5$
$C(2, 3)$	$C(1, 1)$	$C(1, 1)$	$25+1=26$ $4+1=5$
$C(1, 1)$	$C(1, 1)$	$C(1, 1)$	$4+1=5$ $4+4=8$

$$A = [3, 1, 2]$$

$$B = [1, 2, 3]$$

Symmetric distance

$$S_1 \quad d(C_3, C_1) = 16 + 4 = 20$$

$$S_2 \quad d(C_1, C_2) = 1 + 9 = 10$$

$$S_3 \quad d(C_2, C_3) = 0 + 25 = 25$$

$$D(A, B) = 20 + 10 + 25 = 55$$

Clustering - DBSCAN Tutorial

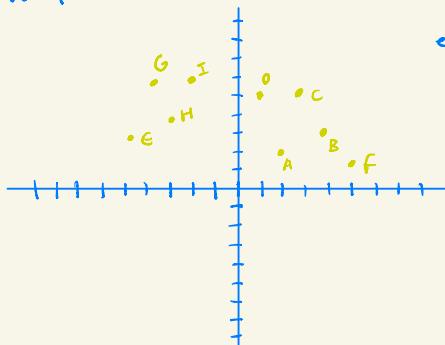
Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

Question 1

Q: Given the data points below and the parameters:
 $\epsilon=3$, $\text{minPoints}=4$, write down which of the points
are core points, border points and noise points,
respectively.

A(2,2) B(4,3) C(3,5) D(1,5)
E(-5,3) F(5,1) G(-4,6) H(-3,4) I(-2,6)

$\epsilon_{PS} = 3$
minpoints < 4



Start from A
A
A is core
 ϵ_{PS} -radius: 3 > minPoints

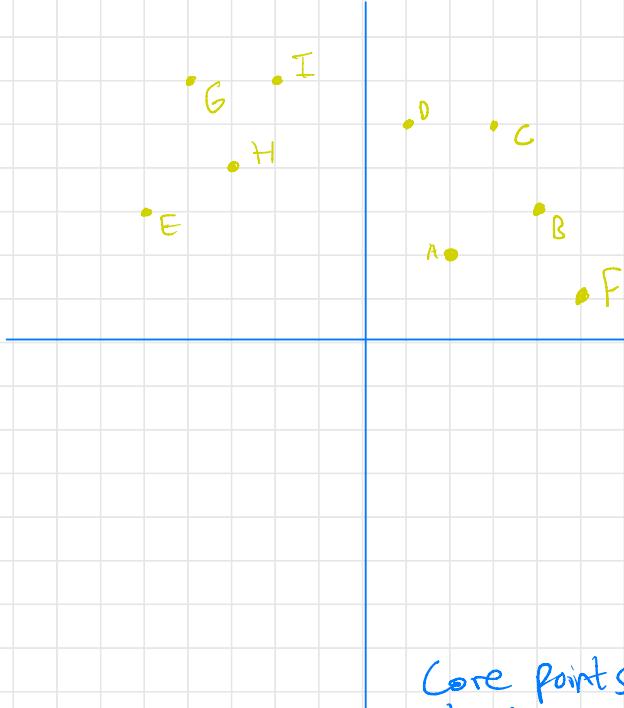
$q \in A$
 $q = BCD(F)$

Solution

DIGSON parameters: epsilon=3, minPoints=4

To simplify calculation, here we use squared L2 distance,
so we need to compare the squared distance with squared epsilon: $\epsilon^2 = 3^2 = 9$
to identify the number of points in the epsilon-neighbourhood:

	A	B	C	D	E	F	G	H	I	#points within epsilon
A	0	5	10	10	50	10	52	29	32	2
B	5	0	5	12	81	5	73	50	45	4
C	10	5	0	4	68	20	50	37	26	3
D	10	23	4	0	40	32	26	17	10	2
E	50	81	68	40	0	104	10	5	18	2
F	10	5	20	32	104	0	106	73	74	2
G	52	73	50	26	10	106	0	5	4	3
H	29	59	37	27	5	73	5	0	5	4
I	32	45	26	10	18	74	4	5	0	3



A is not core point
is Border point
(reach from B)

B is core point
 $\text{eps_rad} : 4 \geq \text{minpoints}(4)$ (B,A,C,F)

F is border point to b

(C,B,D) C is border point $\text{eps_rad} : 3 \leq \text{minpoints}(4)$
(C,D) D is noise point $\text{eps_rad} : 3 \leq \text{minpoints}(4)$
(G,H,I) I is border point $\text{eps_rad} : 3 \leq \text{minpoints}(4)$
(H,I,G,E) H is core point $\text{eps_rad} : 4 \geq \text{minpoints}(4)$
G is border point
E is border point

Core points = B, H
border points = A, C, ~~E~~, F, G, I
Noise points = D

Graph community detection Tutorial

Lin Guosheng
School of Computer Science and Engineering
Nanyang Technological University

Question 1

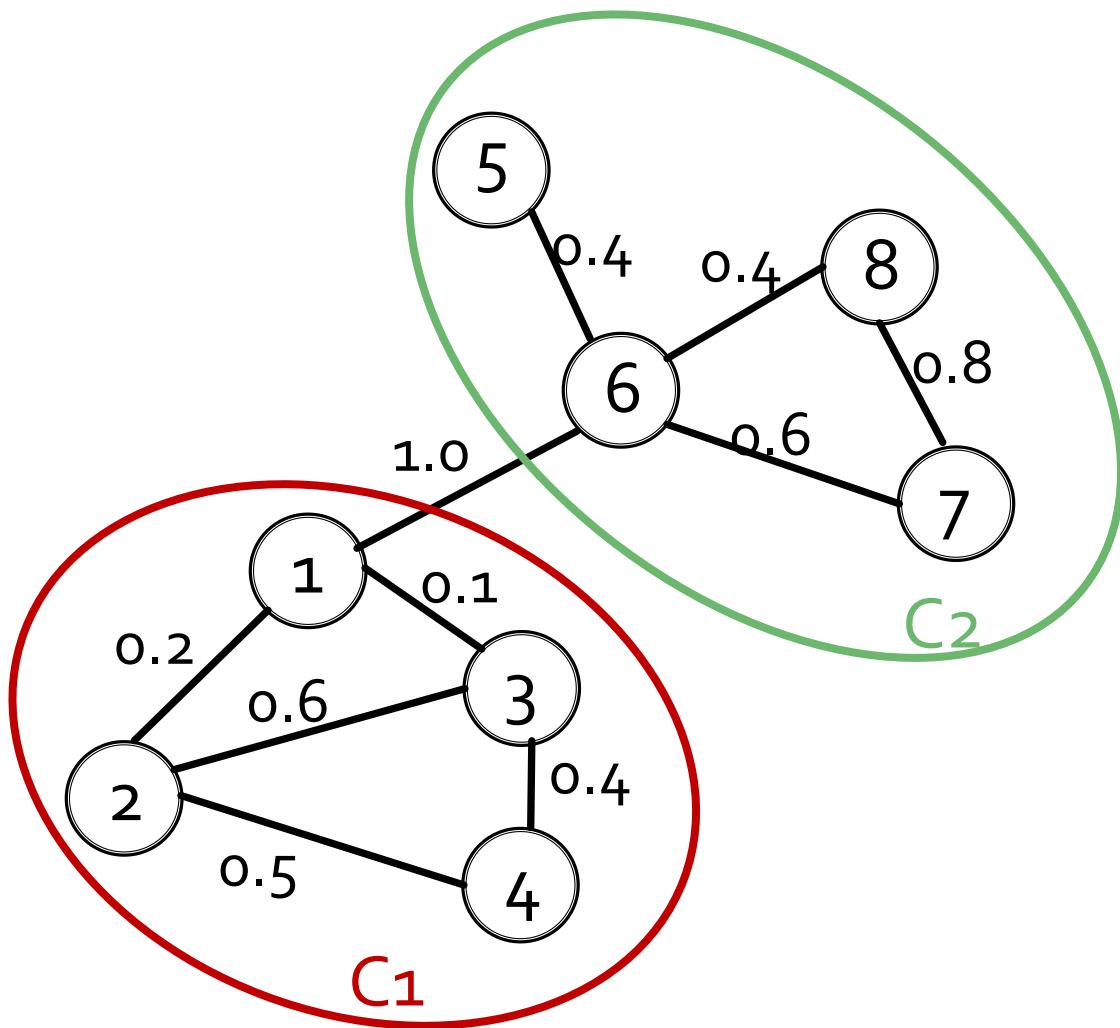
Q: A graph is given below.

The current community assignment is:

$$C1=\{1, 2, 3, 4\},$$

$$C2=\{5,6,7,8\}.$$

Process node #1 to update the communities using Louvain Algorithm.



Process node

$$\Delta Q(1 \rightarrow 2) = Q_{\text{after}} - Q_{\text{before}}$$

$$M = 1.8 + 1 + 2 = 5 \quad = Q(1+2) - [Q(1) + Q(2)]$$

$$Q(1) = -\left(\frac{k_1}{2M}\right)^2 = -\left(\frac{1.3}{2(5)}\right)^2 = -0.0169$$

$$Q(2) = -\left(\frac{k_2}{2M}\right)^2 = -\left(\frac{1.3}{2(5)}\right)^2 = -0.0169$$

$\Sigma_{in} =$

$$0.04 - 0.0676 = -0.0276$$

$$Q(\{1,2\}) = \frac{\Sigma_{in}}{2M} - \left(\frac{\Sigma_{tot}}{2M}\right)^2 = \frac{2(-0.2)}{2.5} - \left(\frac{2(-0.6)}{2.5}\right)^2 =$$



$$\begin{aligned} \Delta Q(1 \rightarrow 2) &= Q_{\text{After}} - Q_{\text{before}} \\ &= -0.0276 - [-0.0169 - 0.0169] \\ &\in 0.0062 \end{aligned}$$

let

affinity matrix $A =$

0 1 0 0	A' =	0 1 0 0
1 0 1 2	0.25 0 0.25 0.5	
0 1 0 0	0 1 0 0	
0 2 0 0	0 1 0 0	

$$\begin{aligned} \Delta Q(1 \rightarrow \{5,6,7,8\}) &= Q_{\text{after}} - Q_{\text{before}} \\ &\in [Q(\{1,5,6,7,8\}) + Q(\{2,3,4\})] - \\ &\quad [Q(\{5,6,7,8\}) + Q(\{1,2,3,4\})] \end{aligned}$$

$$X = \begin{pmatrix} 1 & 0 & -1 & 1 \\ -1 & 1 & 0 & 0 \end{pmatrix}$$

node 2

CZ4032 Data Analytics and Mining

Tutorial for Week 7: association rules

Q1 The following is a pseudocode of Apriori algorithm

- 1) Generate frequent itemsets of length 1
- 2) Repeat until no new frequent itemsets are identified
 - a) Generate length (k+1) candidate itemsets from length k frequent itemsets
 - b) Prune candidate itemsets containing subsets of length k that are infrequent
 - c) Count the support of each candidate by scanning the DB
 - d) Eliminate candidates that are infrequent, leaving only those that are frequent

Explain 2 b) with an example.

Q2 Illustrating Apriori Principle with an example.

Q3 In the lecture, we introduce how to generate length (k+1) candidate itemsets from length k frequent itemsets. Explain this with an example. Can give another way of generating candidates?

Q4 A dataset with 4 records

- 1: a, c,d
- 2, b, c,e
- 3, a, b, c, e
- 4, b, e,

Suppose minimum support is 2. Find frequent itemsets step-by-step.

Q5 Suppose {B,C,D} is a frequent itemset. Enumerate the candidate rules:

Q6 Consider the observation: If $A,B,C \rightarrow D$ is below confidence, so is $A,B \rightarrow C,D$. Can we design an efficient order of generating rules based on the observation?

Go to wooclap.com and use the code **EKDHN0**

How do you generate C2 from F1?

The image shows a laptop screen on the left and a smartphone screen on the right, both displaying a Wooclap poll interface.

Laptop Screen (Left):

- Header: "Go to wooclap.com and use the code **EKDHN0**"
- Question: "How do you generate C2 from F1?"
- Option 1: "Generate all the different pair combinations for all items in F1" (selected, 97% voted)
- Option 2: "Generate all the different pair combinations for all items in F1" (unselected)
- Footer: "wooclap" logo and various icons

Smartphone Screen (Right):

- Header: "pp.wooclap.com" and "wooclap" logo
- Text: "For the previous example, can you design an algorithm that finds frequent itemsets in two passes of the data?"
- Status: "Waiting for next clap"
- Buttons: "Yes" (blue) and "No" (pink)

CZ4032 Data Analytics and Mining

Association Rule Mining: Tutorial 2

Q1 Explain the following observation for PCY algorithm

- a. If a bucket contains a frequent pair, then the bucket is surely frequent
- b. However, even without any frequent pair, a bucket can still be frequent

Q2 Given a dataset, minsup threshold, which of the following has the largest number of itemsets? Which has the smallest number of itemsets?

- 2) Frequent itemsets
- 3) Maximal frequent itemsets
- 4) Closed frequent itemsets

 Q3 Discuss the impact of the following characteristics of a transaction table on the use of the FP tree to mine frequent itemsets from the table:

- (a) Number of unique items in table
- (b) Average number of items in a transaction
- (c) Number of transactions in table

 Q4

A database has four transactions. Let $\text{min_sup} = 60\%$ and $\text{min_conf} = 80\%$.

TID	Date	Items_bought
T100	20006-01-01	{K, A, D, B}
T200	20006-01-01	{D, A, C, E, B}
T300	20006-01-01	{C, A, B, E}
T400	20006-01-01	{B, A, D }

- (a) Find all frequent itemsets using FP-growth.

Q5: for sequence pattern mining, answer the following questions

- Can $\langle\{a\}, \{b\}, \{c\}\rangle$ merge with $\langle\{b\}, \{c\}, \{f\}\rangle$?
- Can $\langle\{a\}, \{b\}, \{c\}\rangle$ merge with $\langle\{b,c\}, \{f\}\rangle$?
- Can $\langle\{a\}, \{b\}, \{c\}\rangle$ merge with $\langle\{b\}, \{c,f\}\rangle$?
- Can $\langle\{a,b\}, \{c\}\rangle$ merge with $\langle\{b\}, \{c,f\}\rangle$?
- Can $\langle\{a,b,c\}\rangle$ merge with $\langle\{b,c,f\}\rangle$?
- Can $\langle\{a\} \{b\} \{a\}\rangle$ merge with $\langle\{b\} \{a\} \{b\}\rangle$?
- Can $\langle\{b\} \{a\} \{b\}\rangle$ merge with $\langle\{a\} \{b\} \{a\}\rangle$?

The PCY algorithm may run slower than the Apriori algorithm at the same minsup setting. Is this statement true?

① Heads up, voting is closed

Yes

No

For step 2 of FP-tree algorithm, if we skip the step "Construct Ordered-Item set", and we only use the alphabetic order of items to construct FP-tree, can we get the correct results (the same set of frequent itemsets)?

① Heads up, voting is closed

Yes

No

Is it always possible to have leaf node containing training data from a ..

① Yes

44% 36

② No

56% 46

wooclap

82 / 103

Will this type of PruneRule affect the generated classifier?

① Yes

86% 84

② No

14% 14

wooclap

98 / 134

Q4 CBA-RG

The CBA-RG Algorithm

```
1   $F_1 = \{\text{large 1-ruleitems}\};$ 
2   $CAR_1 = \text{genRules}(F_1);$ 
3   $prCAR_1 = \text{pruneRules}(CAR_1);$ 
4  for ( $k = 2; F_{k-1} \neq \emptyset; k++$ ) do
5     $C_k = \text{candidateGen}(F_{k-1});$ 
6    for each data case  $d \in D$  do
7       $C_d = \text{ruleSubset}(C_k, d);$ 
8      for each candidate  $c \in C_d$  do
9         $c.\text{condsupCount}++;$ 
10       if  $d.\text{class} = c.\text{class}$  then  $c.\text{rulesupCount}++$ 
11     end
12   end
13    $F_k = \{c \in C_k \mid c.\text{rulesupCount} \geq \text{minsup}\};$ 
14    $CAR_k = \text{genRules}(F_k);$ 
15    $prCAR_k = \text{pruneRules}(CAR_k);$ 
16 end
17  $CARs = \bigcup_i CAR_i;$ 
18  $prCARs = \bigcup_i prCAR_i;$ 
```

C_d is the subset of candidates in C_k that are covered by d

РЧР

S1 2022-2023

2) a) - need to specify the initial

$$b) h_{1a} = 1.0 + 2.1 + 2 - 1 = 0$$

$$h_{1b} = -2.0 + 1.1 + 4.1 - 0 - 3 \rightarrow 0$$

$$h_{2a} = 1.1 + 2.0 + 2.0 = 1 \rightarrow \text{diff from a}$$

$$h_{2b} = -2.1 + 1.0 + 4.0 = -2 \rightarrow 0$$

hamming dist = 1

c)

	Case	C	W	default
r_2	3	3		Y
r_1	3	3		Y
r_5	2	1	1	Y
r_4	3	2	1	N
r_3				

error
 $0+3=3$

$0+3=3$

$1+2=3$

$1+1=2$

A	B	C	Label
2	4	2	N
4	2	1	Y
4	4	2	Y
2	2	1	Y
1	2	1	N
2	3	3	N
4	1	4	N
4	3	3	N
3	4	2	Y
3	1	4	Y
3	3	4	N
3	2	2	N
2	4	1	Y
2	1	1	Y