

ARINJAY-JAIN (A20447307)

Final exam - Data Mining 422.

Problem 1.1

- DBSCAN \rightarrow Define density in a region ~~for~~ feature space based on a center-based methodology
- Density of a point in a region is defined by the number of points with a radius of itself \rightarrow hyperparameter of ~~the~~ radius.

Definition \rightarrow Radius : Eps.

Regions based on Eps.

Dense Region I = Core point
" Edge = Border points
Low / No Density = noise points
(Background)

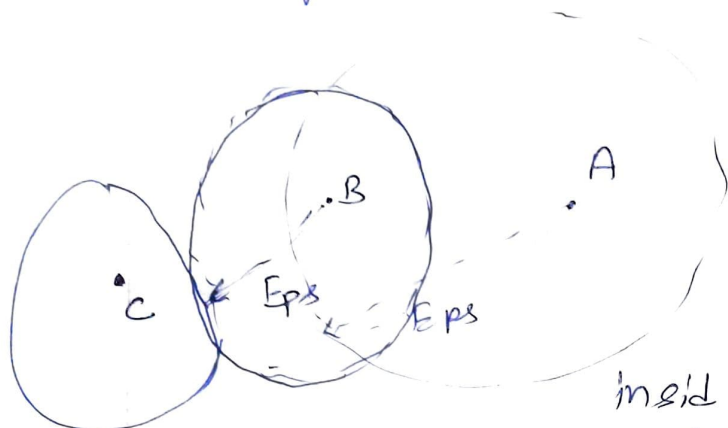
Core points within Eps radius the number of points including it self is $> \text{minpts}$. \rightarrow use define.

Border points : within a neighborhood Eps radius of a core point, but does not have enough points around it to be a core point.

Noise points : Neither enough points with Eps to be a core points and ~~not~~ not ~~se~~ near within Eps to a core points.

$$Eg = \text{Eps} = 10$$

$$\text{min pts} = 7$$



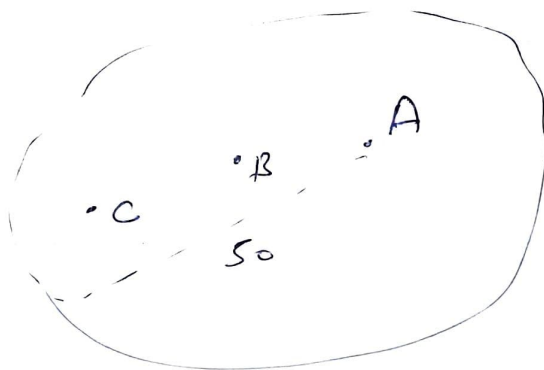
A = core point
B = border point
C = Noise point.

Noise point should not exist inside of the neighborhood a core point. But if we taken a large value of Eps (radius) then it could be possible we capture all the points in one

Cluster.

eg \Rightarrow

$$\text{Eps} = 50$$



1.2)

total points $\Rightarrow n$
features = d

Multivariate Normal

\rightarrow Mahalanobis distance

$$d(x, \bar{x}) = (x - \bar{x}) \Sigma^{-1} (x - \bar{x})^T$$

Covariance

matrix

	x_1	x_2
x_1	σ_1^2	$\sigma_1 \sigma_2$
x_2	$\sigma_1 \sigma_2$	σ_2^2

maximum likelihood estimator:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T \quad (\text{average of } d \times d \text{ matrices.})$$

$$N(x, \mu, \Sigma) = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{-1/2} \exp \left\{ -\frac{1}{2} (x - \mu) \Sigma^{-1} (x - \mu)^T \right\}$$

log likelihood

$$= \log \left(\frac{1}{(2\pi)^{d/2}} \right) + \log |\Sigma|^{-1/2} + \underline{\underline{-\frac{1}{2} (x - \mu) \Sigma^{-1} (x - \mu)^T}}$$

If we remove ~~the~~ an anomaly then log-likelihood may ~~increase~~ decrease ~~or not~~ because we removed some part of distance magnitude. like distance contribution by anomaly will not count in $(x-\mu)^T \Sigma^{-1} (x-\mu)^T$.

If the features are uncorrelated then variance (Correction distance) will ~~in~~ increase so we will get high magnitude.

1.3) Root mean square error (RMSE)

Assume that R is an n by m utility matrix with some entries blank while U and V are matrices of dim n by d and d by m for some d . Also let $P = UV$, the product matrix P .

- 1) Subtract each value of the calculated matrix from the original
- 2) Square each value in the new matrix
- 3) add each row together.
- 4) sum all values the n by 1 matrix
- 5) divide the by the total number of provided rating from our original utility matrix
- 6) take the square root of this result and the value is the ~~RMSE~~ RMSE.

1.4

$$N = 100,000 = 10^6$$

$$\% \text{ of } N = 10^6 \times 10^{-2} = 10^4 = n;$$

length of sing. doc = 1000

$$f_{ij} = 17$$

$$Tf_{ij} = \frac{17}{1000} = \frac{17}{1000} \quad IDF = \log \frac{10^6}{10^4} = \log 10^2 = 2$$

$$Tf-IDf = \frac{17}{1000} \times 2$$

$$= \frac{34}{1000} = 0.034$$

1.6

$$N = 10000000 = 10^7 \quad k = 5 \quad 27 \text{ letters}$$

$$\text{Char matrix dimensionality} = 27 \times 10000000$$

$$\text{Characteristic.} = 27 \times 10^7$$

Signature matrix. ~~5~~ ~~5 \times 10^7~~

~~we have~~ we have 5 permutation.
and 5 shingles.

$$\text{then dim} = 50 \times 5$$

$$r = 5$$

$$b = 10$$

$$P(D_1 \& D_2 \text{ identical in a particular band}) = (0.8)^r$$

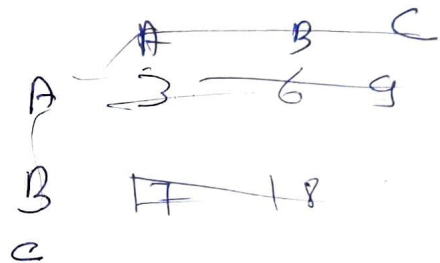
$$= (0.8)^5 = 0.328.$$

1.7

$$A = (3, 6, 9)$$

$$B = (7, 18, 19)$$

$$C = (-2, -4, -8)$$



$$\text{Center of } A = \frac{3+6+9}{3} = \frac{18}{3} = 6$$

$$\text{Center of } B = 18$$

$$\text{center of } C = \frac{14}{3} = 4.66$$

$$\text{Single link} = \min(A, B), (B, C), (A, C)$$

$$\text{single link} = (A, C) = 1.34$$

$$\text{complete} = \max(A, B), (B, C), (A, C)$$

$$= (B, C) = 13.34$$

1.8) The objective of minimizing squared error
causes K-mean to break the large cluster
Hence K-means is not good at deriving
clusters of variable sizes at least when
they are not well separated, hence

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} d(C_i, x)^2$$

also means we are trying to maximize cluster
cohesion.

$$C_k = \frac{1}{n_k} \sum_{x \in C_k} x$$

Optimal centroid value for over optimization solution
is the mean of each cluster!

post processing \rightarrow

Aim is to improve over clustering (reduce SSE)
after the algorithm has completed.

- Split a cluster.
- Introduce a new cluster centroid.

SSE \downarrow if we reduce the K (no. of clusters)
but it may increase the distance b/w
Centroid C_i and x (datapoint) so it also a
cause of \uparrow SSE again.

I. 5)

Dimensionality of Markov matrix = $N \times N$, web pages.
↓

Each col represents the out links from a web page. Hence the sum is equal to 1

Bit Spider trap

is a set of nodes with no dead ends and out links are within in the trap.

$$\text{Take } V_0 = \left(\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right)^T$$

$$M^\infty V_0 = \begin{pmatrix} 0, 0, 1, 0, \dots, 0 \end{pmatrix}^T$$

Hence everyone ends up at K nodes.
for K node page rank score is 1.

for $(n-K)$ node the page Rank score will be Zero (0).

to avoid this problem we can method called taxation.

2.1)

Silhouette score / coeff.

for a set of sample data points is used to measure how dense and well separated the clusters are.

mean intra-cluster distance (a) : the avg distance from x_i to other points in C

mean inter(nearest)-cluster distance (b) : the avg distance from x_i to other points in D .

Silhouette score

$$S = \frac{b_i - a_i}{\max(b_i, a_i)}$$

$$\text{range} = [-1, 1]$$

lowest value (-1) tells us the samples are assigned in the wrong clusters,

high value (1) : cluster is dense and well-separated
we want $a_i \rightarrow 0$ to get S_i to 1 .

zero (0) \Rightarrow overlapping cluster with sample very close to the decision boundary of the neighboring cluster.

2.2)

$$\mu = 100$$

$$\sigma^2 = 10000 \quad SD = \sqrt{10^4} = 100$$

$$\Pr(\mu - 4\sigma \leq X \leq \mu + 3\sigma) \approx 0.01\% \\ \approx 0.0001$$

$$Z \text{ score} = Z = \frac{X - \mu}{\sigma}$$

↑
attribute value.

$$C.I. = \bar{X} \pm 4 \frac{\sigma}{\sqrt{n}}$$

Confided interval of 4-standard deviation.

give the value above the mean and below the mean.

2.3)

$$\begin{array}{rcccccc} U_1 = & 4 & 2 & 3 & 2 & 4 \\ U_2 = & 5 & 3 & 4 & 3 & 5 \end{array}$$

$$\text{Avg } U_1 = \frac{15}{5} = 3$$

$$\text{Avg } U_2 = \frac{20}{5} = 4$$

Normalize centre \Rightarrow

$$\begin{array}{rcccccc} U_1 = & 1 & -1 & 0 & -1 & 1 \\ U_2 = & 1 & -1 & 0 & -1 & 1 \end{array}$$

okay both user profile is similar.

$$\text{Cos: Similarity} = \text{Corr} = \frac{(1 \times 1) + (-1 \times -1) + (0) + (-1 \times -1) + (1 \times 1)}{\sqrt{1^2 + (-1)^2 + (0)^2 + (-1)^2 + (1)^2} \times \sqrt{1^2 + (-1)^2 + (0)^2 + (-1)^2 + (1)^2}}$$

$$= \frac{1 + 1 + 1 + 1}{2 \times \sqrt{4}}$$

$$= \frac{1 + 1 + 1 + 1}{4} = 1$$

$1 > 0.75 \Rightarrow$ highly correlated.
 \uparrow
 threshold.

2.4) In Content-Based-Filtrering.

Recommnd item_i to user U based upon other item_j.

$i_1, i_2 \dots i_k$ rated by user U.

methodology \Rightarrow

- Build Item profile based on Item properties
 \rightarrow we see what a User likes based on Item profiles of brought / rated item \rightarrow this gives us a user profile

\rightarrow we find Item profile that match user profile for recommendation.

User profile utility matrix.

eg \Rightarrow family animation adventure drama docum.

U_1

	M_1	M_2	\dots	M_{20}	\dots	M_{100}	
U_1	3	5		2		7	\leftarrow any rating

Item profile

family

anima

adu

drama

doc.

genre
genre

i_1

0

1

0

0

0

i_2

1

0

0

0

0

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

\vdots

i_{100}

\vdots

\vdots

\vdots

\vdots

\vdots

\rightarrow movie

unweighted

	m_1	m_2	m_3	...	m_n
U_1	3	4	5		

Average for

user $u = 3$

for non doc genre movie.

weighted for (non-doc movie) ← movie

	m_1	m_2	m_3	...
user → U	3-3	4-3	5-3	...
	0	1	2	...
	100			

Average for

user $u = 5$

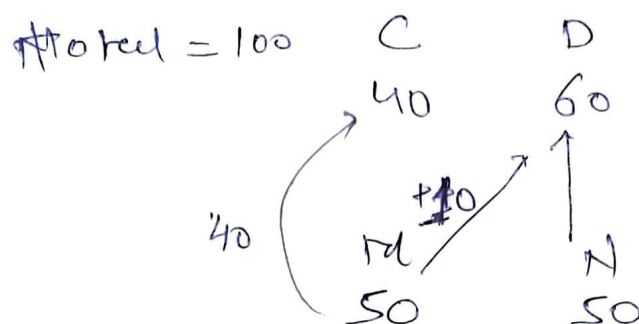
for doc genre movie.

	m_1	m_2	m_3	...
user → U	3-5	4-5	5-5	...
	-2	-1	0	...

$$= -3/3 = -1$$

$$= \frac{-3 + 0 \dots}{100}$$

3.1)



$$E = - \sum P(w_c) \log_2 P(w_c)$$

$$= - \sum \frac{|w_c|}{nw} \log_2 \frac{|w_c|}{nw}$$

$|w_c|$ = is the count of points classified as c in cluster w .

nw = is the count of points in cluster w .

Entropy D $\Rightarrow nw = 60$

$$|w_{c1}| = 50 \text{ (class N)}$$

$$|w_{c2}| = 10 \text{ (class M)}$$

$$\# \text{ Entropy D} = - \frac{50}{60} \log_2 \frac{50}{60} - \frac{10}{60} \log_2 \frac{10}{60}$$

Similarly

$$\# \text{ Entropy C} = - \frac{40}{40} \log_2 \frac{40}{40} - \frac{0}{40} \log_2 \frac{0}{40}$$

$$= 0$$

Q Lucky 7

- 1) 2018 ACM AM Turing Award
- 2) Recipes for pizza
- 3) \$432500
- 4) openAI
- 5) Pong
- 6) ~~Hanabi~~ Hanabi
- 7) Finland