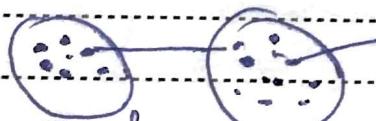
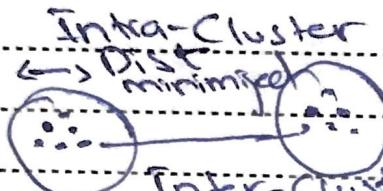


DATE

DAY

MON TUE WED THU FRI SAT SUN

Cluster: Collection of data Objects



familiar
within
group

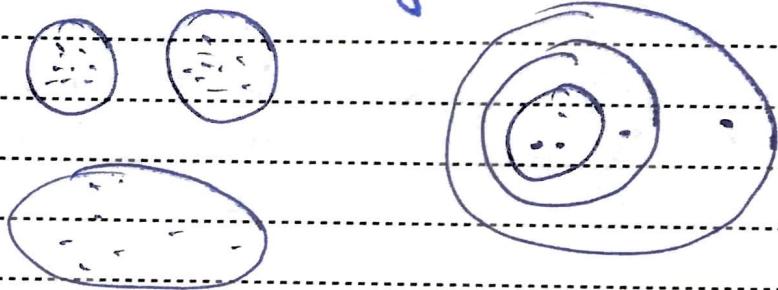
Inter-Cluster
Dist maximized

Dissimilar
outside group

Original Points

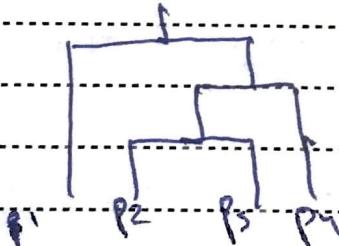
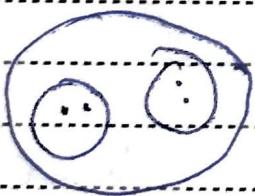
Partitional Clustering

Traditional
Hierarchical



Non-Traditional
Hierarchical

Traditional
Dendogram



Partitioning Method: Partitioning a database of n objects into a set of k clusters, such that the sum of sq dist is minimized where c_i is the centroid or medroid.

$$E = \sum_{i=1}^k \sum_{p \in C_i} (p - c_i)^2$$

 stewart

DATE

DAY

MON TUE WED THU FRI SAT SUN

Heuristic Method: K-means, K-medoids

K-Means Clustering

$O(n + k \times I \times d)$ n = no. of points, k = no. of clusters, I = no. of iterations, d = no. of attributes

	Data Points	D: A1	D: B1	D: C1	C
A1	(2, 10)	0	3.61	8.06	1
	2, 5	5	4.24	3.16	3
B1	8, 4	8.5	5	7.28	2
	5, 8	3.61	0	7.21	2
C1	7, 5	7.07	3.61	6.71	2
	6, 4	7.21	4.12	5.39	2
A1	1, 2	8.06	7.21	0	3
	4, 9	2.24	1.41	7.62	2

Step 1: Select Random Centroids

Step 2 Find Dist of each Point for every cluster

Step 3 Assign the point to the closest cluster

Step 4 Update the centroids

Step 5 Do until converge

Euclidean

$$\text{Dist} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

DATE

&

DAY

MON TUE WED THU FRI SAT SUN

New Centroids

$$A_1 = (2, 10)$$

$$B_1 = \left(\frac{8+5+7+6+4}{5}, \frac{4+8+5+4+9}{5} \right) = (5, 8)$$

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

Data Points	D:A1	D:B2	D:C1	C
2 10	0	5.66	6.52	1
2 5	5	4.12	1.58	3
8 4	8.49	2.83	6.52	2
5 8	3.61	2.24	5.70	2
7 5	7.07	1.41	5.70	2
6 4	2.21	2	4.53	2
1 2	8.06	6.4	1.58	3
4 9	2.24	3.61	6.04	1

Do until all clusters converge

K-Node Clustering

For categorical data

DATE

DAY

MON TUE WED THU FRI SAT SUN

we find the ~~disse~~ dissimilarities b/w
the point and the cluster

	Person	hair color	eye color	skin color	P:A1	O:B1	O:C1
A1	P1	blonde	amber	fair	0	2	2
	P2	brunette	gray	brown	3	3	3
	P3	red	green	brown	3	1	3
	P4	black	hazel	brown	3	3	1
	P5	brunette	amber	fair	1	2	2
2	P6	black	gray	brown	3	3	2
	P7	red	green	fair	2	0	2
C1	P8	black	hazel	fair	2	2	0

New Centroids

A1	Blonde	amber	fair	1
	brunette	gray	brown	2
	brunette	amber	fair	3

Since blonde occurs 2 time we will select blonde. Same with amber and fair. If all were diff then we would've selected randomly.

DATE

Day

MON TUE WED THU FRI SAT SUN

A1 : blonde amber fair

Similarly do with other points

B1 : red green fair

C1 : black hazel brown

Repeat the process until Converge

k-Mediod

X	Y	D: C1	D: C2	Cluster
8	7	6	2	2
3	7	3	7	1
4	9	4	8	1
9	6	6	2	2
8	5	-	-	
5	8	4	6	1
7	3	5	3	2
8	4	5	1	2
7	5	3	1	2
4	5	-	-	

$$\text{Dist} = |x_1 - x_2| + |y_1 - y_2|$$

DATE

&

DAY

MON TUE WED THU FRI SAT SUN

$$\text{cost} = (3+4+4) + (3+1+1+2+2) = 20$$

X	Y	D:C1	D:C2	Cluster
8	7	6	3	2
3	7	3	8	1
4	9	4	9	1
9	6	6	3	2
8	5	4	1	2
5	8	4	7	1
7	3	5	2	2
C1	8	4	-	-
7	5	3	2	2
C2	4	5	-	-

New

$$\text{cost} = (3+4+4) + (2+2+1+3+3) = 22$$

$$\begin{aligned}\text{Swap cost} &= \text{New cost} - \text{Old cost} \\ &= 22 - 20 = 2 > 0\end{aligned}$$

As the swap cost is not less than zero, we undo the swap

DATE & DAY

MON TUE WED THU FRI SAT SUN

Hierarchical Clustering

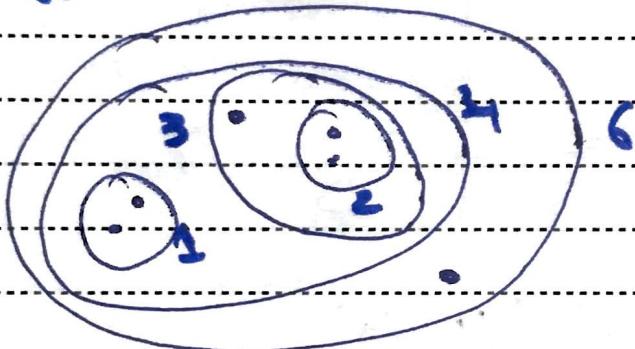
uses dist matrix as a clustering criteria

Hierarchical
Clustering

Agglomerative

Divisive

Agglomerative Clustering



bottom-up
Approach

Consider the following points

18, 22, 25, 42, 27, 43

Step: 1

using
Single
Linkage

	18	22	25	27	42	43
18	0					
22	9	0				
25	7	3	0			
27	9	5	2	0		
42	24	20	17	15		
43	25	21	18	16	1	0

DATE

DAY

MON TUE WED THU FRI SAT SUN

	18	22	25	27	42,43
18	0				
22	4	0			
25	7	3	0		
27	9	5	2	0	
42,43	24	20	17	15	0

	18	22	25,27	42,43
18	0			
22	4	0		
25,27	7	3	0	
42,43	24	20	15	0

	18	25,27,22	42,43
18	0		
25,27,22	9	0	
42,43	24	15	0

	18,25,27,22	42,43
18,25,27,22	0	
42,43	15	0

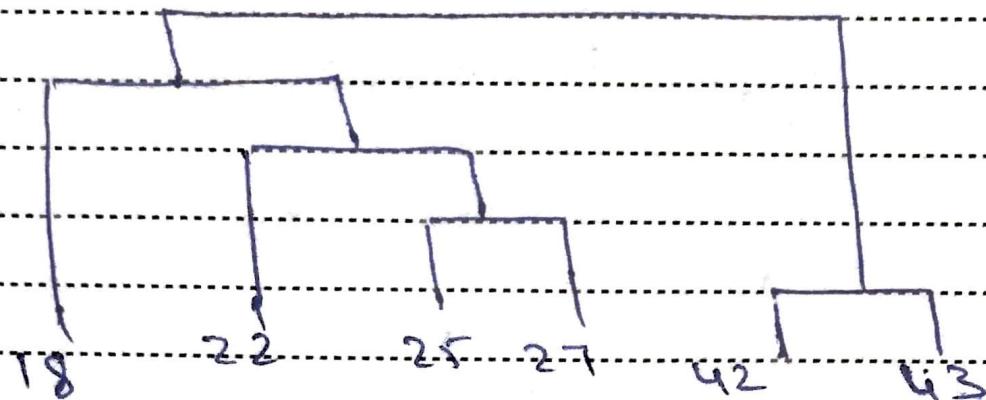
18, 25, 27, 22, 42, 43

DATE

&

DAY

MON TUE WED THU FRI SAT SUN



Single Linkage Min

Complete Linkage Max

Average Linkage Avg of Dist

Divisive Clustering

Top down
Approach

P * Y

Step 1

dist matrix

	1	2	3	4	5	6	
1	0.4	0.53					
2	0.22	0.38					
3	0.35	0.32					
4	0.26	0.19	1	0			
5	0.08	0.41	2	0.23	0		
6	0.45	0.31	3	0.22	0.15	0	
			4	0.37	0.20	0.15	0
			5	0.34	0.14	0.28	0.29
			6	0.23	0.25	0.11	0.22

DATE

&

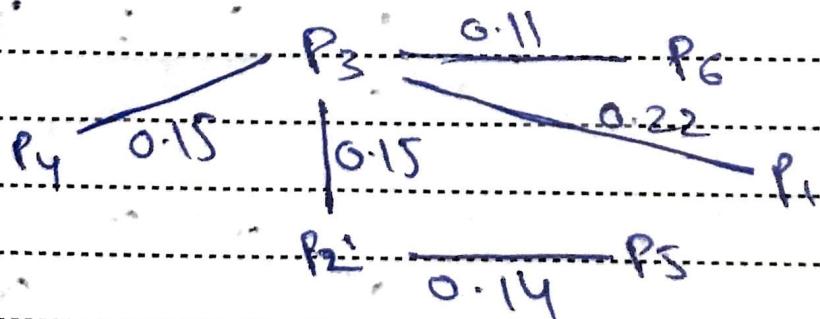
DAY

MON TUE WED THU FRI SAT SUN

Step:2 Compute Minimum Spanning Tree

Edge	cost ← Prim's Algo
P3, P6	0.11
P2, P5	0.14
P2, P3	0.15
P3, P4	0.15
P2, P4	0.20
P1, P3	0.22
P1, P2	0.23
P1, P6	0.23
P2, P6	0.25
P3, P5	0.28
P4, P5	0.29
P1, P5	0.34
P1, P4	0.37
P5, P6	0.39

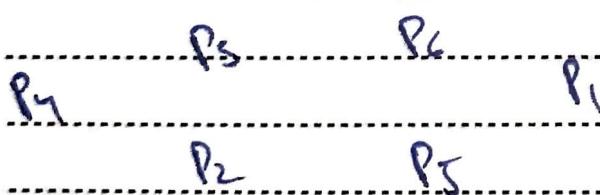
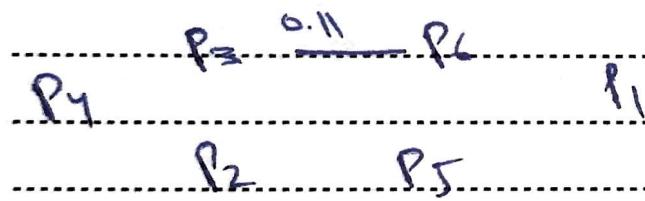
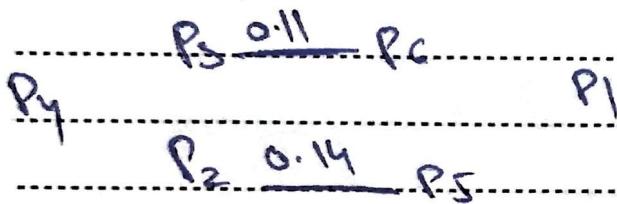
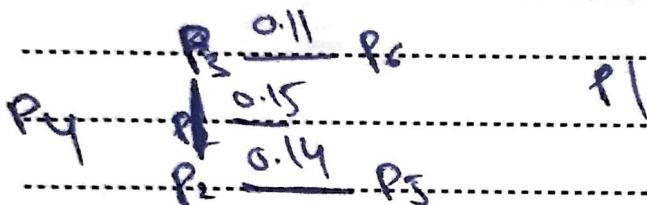
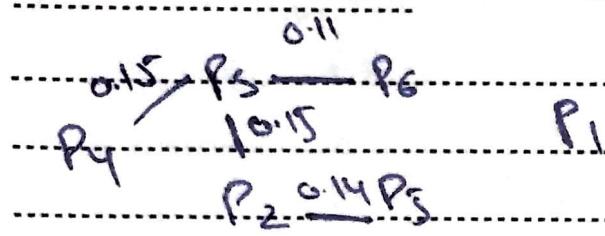
Step:3



DATE

Applying Complete
Linkage to break
the edges

MON TUE WED THU FRI SAT SUN

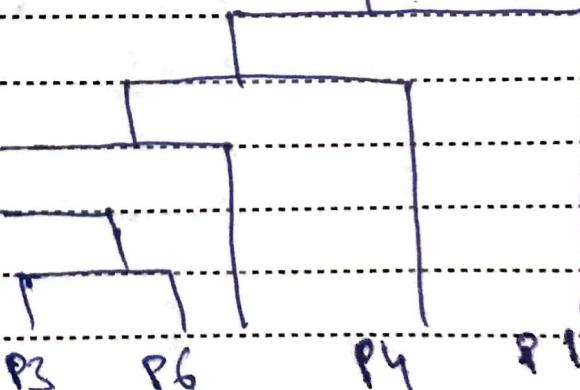


TOP



Down

 stewart™
P2 P5



DATE

23

DAY

MON TUE WED THU FRI SAT SUN

BIRCH

Branching factor $B=2$

max no. of sub-cluster at each leaf node $L=5$

threshold $T=1.5$

Data

x	y	x^2	y^2
3	4	9	16
2	6	4	36
4	5	16	25
4	7	16	49
3	8	9	64
6	2	36	4
7	2	49	4
7	4	49	16
8	4	64	16
8	5	64	25

→ Uses CF (Clustering Feature) to summarize a cluster. CF is a 3D vector summarizing info about cluster objects

$$CF = (n, LS, SS)$$

DATE



DAY

MON TUE WED THU FRI SAT SUN

where n is the no. of objects in the cluster, LS is the linear sum of the objects and SS is the sq. sum of obj.

$$LS = \sum_{i=1}^n x_i$$

$$SS = \sum_{i=1}^n x_i^2$$

cluster centroid $x_0 = LS/n$

cluster Radius $R = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - x_0)^2}$

$$\sqrt{\frac{SS}{n} - \frac{LS^2}{n^2}} \text{ OR } \sqrt{\frac{n(SS) - 2LS^2 - n(LS)}{n^2}}$$

cluster Diameter $D = \sqrt{\frac{\sum_{i=1}^n \sum_{j=1}^n (x_i - w_j)^2}{n(n-1)}}$

$$= \sqrt{\frac{2n(SS) - 2(LS)^2}{n(n-1)}}$$

DATE



DAY

MON TUE WED THU FRI SAT SUN

Consider data point $u_1 = (3, 4)$

$$\text{Radius} = 0$$

$$LS = 3, 4$$

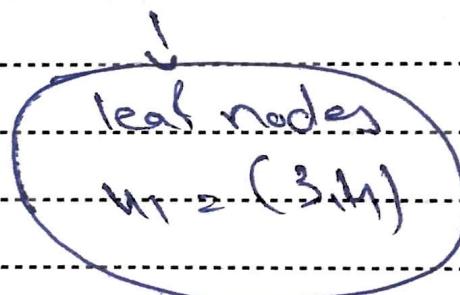
$$SS = 9, 16$$

$$CF = \{1, (3, 4), (9, 4)\}$$

CF: $\{(1, 3, 4), (9, 4)\}$

For data point

$$u_2 = (2, 6)$$



$$LS = (3, 4) + (2, 6)$$

$$= (5, 10)$$

$$SS = (9, 16) + (4, 36)$$

$$= (13, 52)$$

$$\alpha = 2$$

~~Radius~~ $= \sqrt{\frac{13,52}{2} - (5,10)^2}$

$$= \sqrt{(6.5, 26) - (6.25, 25)}$$

$$= \sqrt{(0.25, 1)} = (0.5, 1)$$

DATE &

DAY

MON TUE WED THU FRI SAT SUN

$$R = (0.5, 1) \subset (T, T) : T = 1.5$$

Therefore u_2 will cluster with leaf u_1 .

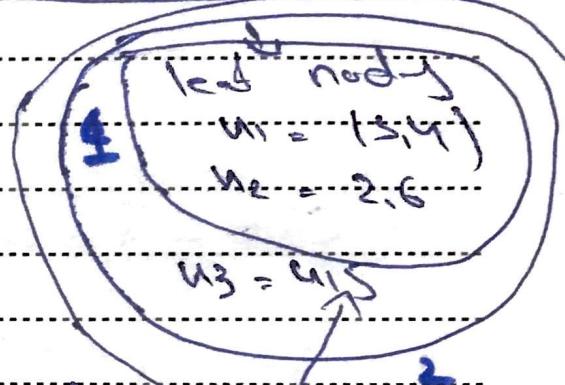
$$CFI \subset 2, (5, 10), (13, 52) >$$

For $u_3 = (u_1)$

$$CFI \subset 2, (5, 10), (13, 52)$$

$$\begin{aligned} LS &= (4, 5) + (5, 10) \\ &= (9, 15) \end{aligned}$$

$$\begin{aligned} SS &= (16, 25) + (13, 52) \\ &= (29, 77) \end{aligned}$$



$$R = (0.47, 0.47) \subset (1.5, 1.5) \quad T \quad u_1 = 4, 7$$

$$CFI \subset 3, (9, 15), (29, 77) >$$

For $u_4 = (4, 7)$

$$LS = (4, 7) + (9, 15) = (13, 22)$$

$$SS = (16, 49) + (29, 77) = (45, 126)$$

$$R = (0.47, 0.47) \subset (1.5, 1.5) \quad T$$

$$CFI \subset 4, (13, 22), (45, 126) >$$

DATE

DAY

MON TUE WED THU FRI SAT SUN

For $u_5 = (3, 8)$

$$LS = (13, 22) + (3, 8) = (16, 30)$$

$$SS = (9+45, 64+126) = (54, 190)$$

$$R = (0.33, 0.63) \leftarrow (1.5, 1.5) \uparrow T$$

Add the point to
the leaf node

$$CF1 \leftarrow 5, (16, 30), (54, 190) \uparrow$$

for $u_6 = (6, 2)$

$$LS = (\cancel{13}, 30) + (6, 2) = \cancel{13}(22, 32)$$

$$SS = \cancel{13} (54, 190) + (36, 4) = (90, 194)$$

$$R = (1.24, 1.97) \leftarrow (1.5, 1.5) \uparrow F$$

Create a New Leaf

$$CF2 \leftarrow 1, (6, 2), (36, 4)$$

Since $B=2$, only 2 branches will exist

DATE



DAY

MON TUE WED THU FRI SAT SUN

For $u_7 = (7,2)$

$[CF1 \leftarrow S, (16,30), (54,190) \rightarrow CF2 \dots]$

firstly we will check u_7 is close to which branch

leaf nodes

$$u_1 = (3,4)$$

$$u_2 = (2,6)$$

$$u_3 = (4,5)$$

$$u_4 = (4,7)$$

$$u_5 = (3,8)$$

leaf

nodes

$$u_6 = (6,2)$$

$$CF1 = LS/N = ((16,30)/5 \\ = (3.2,6)$$

$$CF2 = (6,2)/1 \\ = 6,2$$

u_7 is close to $CF2$

$$LS = (7,2) + (6,2) = (13,4)$$

$$SS = (49,4) + (36,4) = (85,8)$$

$$R = (0.5, 0) \subset (1.5, 1.5) T$$

$$CF2 = L2, (13,4), (85,8)$$

Add point to $CF2$ leaf node

DATE



DAY

MON TUE WED THU FRI SAT SUN

For $u_8 = (7, 4)$

$$CF_1 = (16, 30)_{15} \\ = (32, 6)$$

$$CF_2 = \frac{(13, 4)}{2}$$

$$= (6.5, 2)$$

Closer to CF_2

$$LS = (13, 4) + (7, 4) = (20, 8)$$

$$SS = (44, 16) + (85, 8) = (134, 24)$$

$$R = (0.47, 0.94) \angle (1.5, 1.5)^\circ$$

$$CF_2 \angle 3, (20, 8), (134, 24) >$$

For $u_9 = (8, 4)$

$$CF_1 = (16, 30)_{15} \\ = (32, 6)$$

$$CF_2 = (20, 8)_{13} \\ = (6.6, 2.6)$$

Closer to CF_2

$$LS = (20, 8) + (8, 4) = (28, 12)$$

$$SS = (64 + 134) + (16 + 24) = (198, 40)$$

 **stewart**

$$R = (0.7, 1) \angle (1.5, 1.5)^\circ$$

DATE & DAY

MON TUE WED THU FRI SAT SUN

$$110 \div (7,9)$$

$$\begin{aligned} CF_1 &= (16, 30) / 5 \\ &= (3, 2, 6) \end{aligned}$$

$$\begin{aligned} CF_2 &= (28, 12) / 4 \\ &= (7, 3) \end{aligned}$$

Closer to CF2

$$LS = (16, 30) + (7, 9) = (23, 39)$$

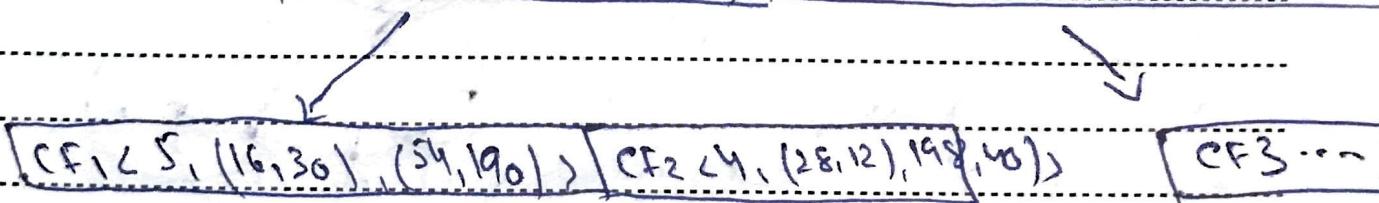
$$SS = (54, 190) + (49, 81) = (103, 271)$$

$$R = (1.57, 1.7) L (1.5, 1.5) F$$

Create a new CF3

$$CF_3 \leftarrow (1, (7, 9), (49, 81)) \rightarrow$$

$$CF_{12} \leftarrow (9, (44, 42), (252, 230)) \quad CF_3 \leftarrow (1, (7, 9), (49, 81)) \rightarrow$$



leaf nodes

$u_1 \dots u_5$

leaf nodes

$u_6 \dots u_9$

leaf nodes

u_{10}

DATE

&

DAY

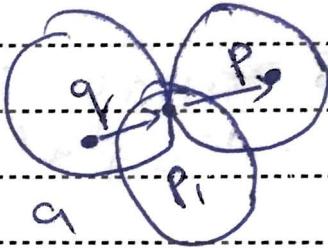
MON TUE WED THU FRI SAT SUN

Density Based Clustering

Eps : Max radius of neighbourhood

Minpts : Min no. of points in an Eps -neighbourhood

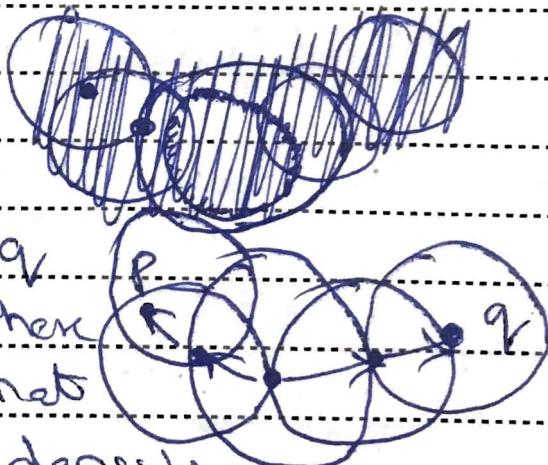
Density-Reachable



A point p is density-reachable from q wrt $\text{Eps}, \text{Minpts}$ if there is a

chain of points p_1, \dots, p_n , $p_1 = q$, $p_n = p$
such that p_{i+1} is directly-density-reachable from p_i

Density-Connected



A point p is density-connected to a point q wrt $\text{Eps}, \text{Minpts}$ if there

is a point o such that both p and q are density-reachable from o wrt Eps and Minpts

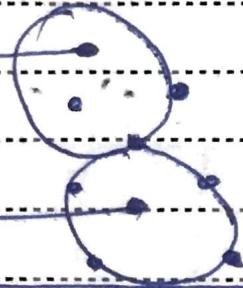
DATE

&

DAY

MON TUE WED THU FRI SAT SUN

Border
point



Outer

Core
point

$\text{EPS} = 1\text{cm}$
 $\text{Minpt} = 5$

Core Point: Obj or point which has atleast
Minpt obj within the radius of ϵ

Boundary Point: Direct Density Reachable.

An Obj or point which is not a core point
but it is in the neighbourhood of the core
point with the radius of ϵ

Noise Point: An Obj/Point which is neither a
core point or boundary point

S1 S7

$$\epsilon = 3.5$$

S2 S4 S9

$$\text{Minpt} = 3$$

S3 3 3

S4 4 4

S5 3 7

S6 6 7

S7 6 1

S8 5 5

DATE



DAY

MON TUE WED THU FRI SAT SUN

Step 1 Distance Matrix

	S1	S2	S3	S4	S5	S6	S7	S8
S1	0							
S2	4.24	0						
S3	4.41	S1	0					
S4	3.16	4	1.41	0				
S5	2	5.83	4	3.16	0			
S6	1	3.61	5	3.61	3	0		
S7	6.08	3.61	3.61	3.61	6.71	6	0	
S8	2	3.16	2.83	1.41	2.83	2.24	4.12	0

Step 2 Identify the neighbours of each point for $\epsilon = 3.5$

S1 : S4, S5, S6, S8

S2 : S8

S3 : S4, S8

S4 : S5, S8, S1, S3

S5 : S1, S4, S6, S8

S6 : S1, S5, S8

S7 : None

S8 : S1, S2, S3, S4, S5, S6

DATE

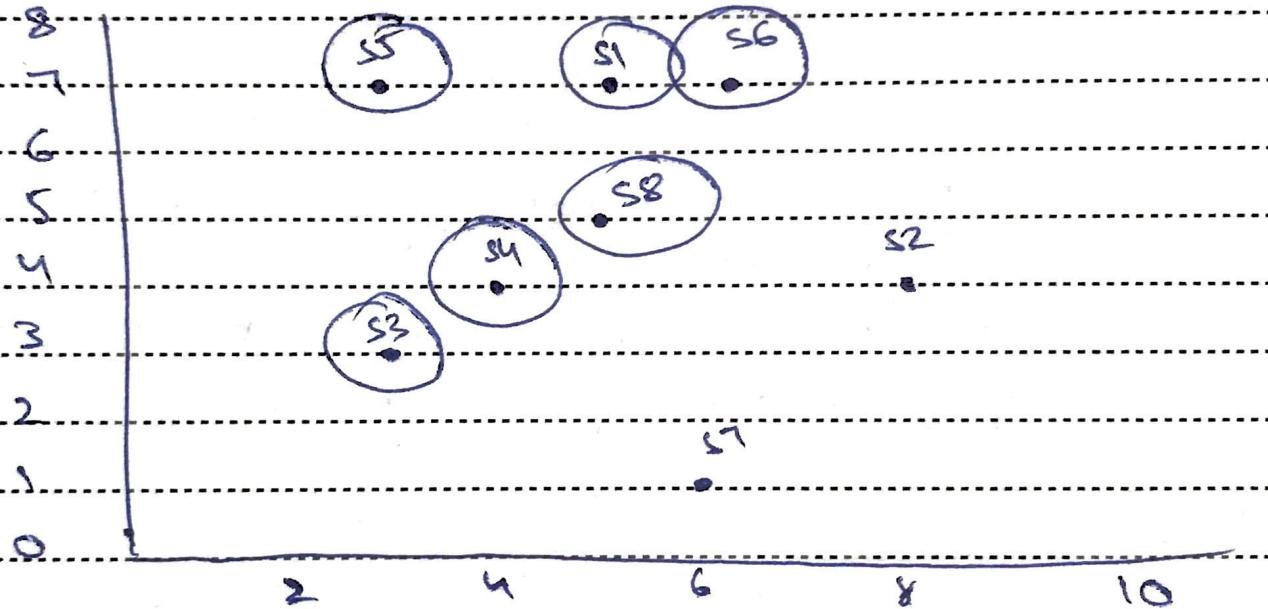
&

DAY

MON TUE WED THU FRI SAT SUN

Step: 3 Ascertain whether the point is
the core point for Minpt = 3

s_1, s_3, s_4, s_5, s_6 and s_8 are core
points



s_2 and s_7 are noise points. Moreover
 s_2 can be classified as a boundary point
since it is direct density reachable
from s_8

DATE

&

DAY

MON TUE WED THU FRI SAT SUN

Text Classification

Bag of words

- Represents text as a collection of words
- Counting occurrences in a document (corpus)

Cons

- Ignore content and meaning of words which causes ambiguity
- cannot capture the b/w word such as synonym etc
- creates sparse feature space which causes overfitting

Example

S1 : The cat in the hat

S2 : The dog chased the cat

S3 : The hat was lost

unique words: The, cat, in, hat, dog, chased, was, lost

vector S1 : 11110000

" S2 : 11001100

" S3 : 10010011

DATE



DAY

MON TUE WED THU FRI SAT SUN

Vector Space Model

→ Represents texts as vectors of features, where each feature represents a term or word in the text and each vector represents a document

Cons

- Sparse Representation of the data
- doesn't capture the context and meaning of word
- suffer from curse of dimensionality

Example

D1: The cat in the hat vector D1: 11110000

D2: The dog chased the cat vector D2: 11001100

D3: The hat was lost vector D3: 10010011

unique words: The, cat, in, hat, dog, chased, was, lost

Similarity b/w

D1 and D2

$$\frac{1*1 + 1*1 + 1*0 + 1*0 + 0*1 + 0*1 + 0*0 + 0*0}{\sqrt{1^2+1^2+1^2+1^2} \sqrt{1^2+1^2+1^2+1^2}}$$

$$= \frac{2}{\sqrt{4*4}} =$$

DATE

&

DAY

MON TUE WED THU FRI SAT SUN

Term Frequency

- How often a term appears in a doc. TF weighting is the process of normalizing the TF by taking into account the length of doc, so that longer docs do not have an unfair advantage over shorter docs.
- Give more imp. to terms that are relevant

$$1 + \log(TF) \text{ OR } \log(TF)$$

Example

	disease	symptom	osteoporosis	
Doc 1	20	15	1	36
Doc 2	3	6	13	26

$$1 + \log(TF)$$

	Disease	symptom	osteoporosis	
Doc 1	2.3	2.2	1	5.5
Doc 2	1.47	2	2.1	5.57

DATE

&

DAY

MON TUE WED THU FRI SAT SUN

Inverse Document frequency

- Measure how imp a term is in the corpus or collection of documents
- IDF weighting gives more weight to terms that are rare in the corpus and less weight to the terms that are common

$$IDF = \log(N/df(t)) \text{ or } 1 + \log(N/df(t))$$

where N is the total no. of docs in collect
and $df(t)$ is the no. of docs containing the
term t

Example

Suppose $N = 100,000$

	$df(t)$	IDF
disease	2000	$1 + \log(100,000 / 2000) = 2.7$
sympton	300	$1 + \log(100,000 / 300) = 3.5$
osteoporosis	10	$1 + \log(100,000 / 10) = 5$

DATE



DAY

MON TUE WED THU FRI SAT SUN

TF-IDF

- Measures the imp of a term in a doc within a corpus
- Gives higher weigh to terms that are frequent in the doc but rare in corpus and lower weight to terms that are frequent in both doc and corpus
- Biased towards longer doc and give more weights to words like "the" or "and"

$$\text{TF-IDF} = \text{TF} * \text{IDF}$$

Example

$$1 + \log(\text{TF})$$

	Disease	Symptom	Osteoporosis
Doc 1	2.3	2.2	1
Doc 2	1.47	2	2.1

	d(TF)	IDF
Disease	2000	2.7
Symptom	300	3.5
Osteoporosis	10	5

$$\text{TF-IDF Doc 1: } 2.3 * 2.7 + 2.2 * 3.5 + 1 * 5 = 18.9$$

$$\text{TF-IDF Doc 2: } 1.47 * 2.7 + 2 * 3.5 + 2.1 * 5 = 21.5$$

DATE



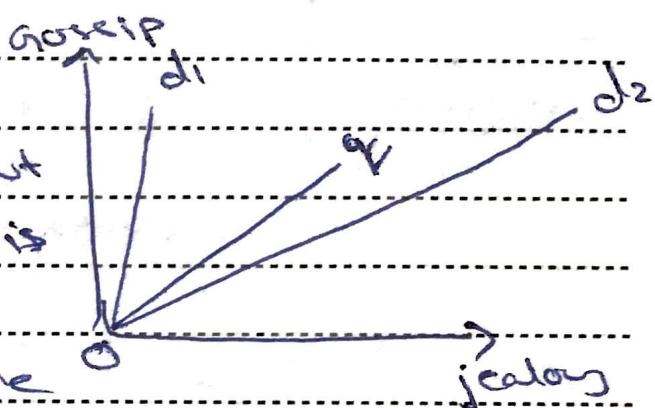
DAY

MON TUE WED THU FRI SAT SUN

Cosine Similarity

- Unaffected by the length of doc
- Robust measure of similarity that can handle high-dimensional vector spaces

Dist of terms in doc 2 is similar to the query q_1 but if dist was used then d_2 is very different from q_1 . Therefore we normalize the length



Angle to cosine

Text is converted into vectors and vector contains count of words. word count cannot be neg. therefore angle to cosine would always be positive and would range from 0-1

$$\cos(\vec{v}_q, \vec{v}_d) = \vec{v}_q \cdot \vec{v}_d$$

$$\cos(\vec{q}, \vec{d}) = \frac{\vec{q} \cdot \vec{d}}{\|\vec{q}\| \|\vec{d}\|}$$

DATE



DAY

MON TUE WED THU FRI SAT SUN

Example 1 Doc to Doc

term	Doc1	Doc2	Doc3
affection	115	58	20
jealous	10	7	11
gossip	2	0	6
wuthering	0	0	38

$1 + \log(TF)$

term	Doc1	Doc2	Doc3
affection	3.06	2.76	2.30
jealous	2.00	1.85	2.04
gossip	1.30	0	1.78
wuthering	0	0	2.58

$$\text{length of Doc1} = \sqrt{(3.06)^2 + (2)^2 + (1.30)^2 + 0^2} = 3.89$$

$$\text{length of Doc2} = \sqrt{(2.76)^2 + (1.85)^2 + 0^2 + 0^2} = 3.32$$

$$\text{length of Doc3} = \sqrt{(2.30)^2 + (2.04)^2 + (1.78)^2 + (2.58)^2} = 6.40$$

length Normalization

DATE



DAY

MON TUE WED THU FRI SAT SUN

Term	Doc1	Doc2	Doc3
affection	$3/3.89 = 0.78$	$276/332 = 0.832$	$230/440 = 0.524$
jealousy	$2/3.89 = 0.51$	0.555	0.465
gossip	$1.30/3.89 = 0.35$	0	0.405
wuthering	0	0	0.588

$$\cos(\text{Doc1}, \text{Doc2}) = 0.78 + 0.832 + 0.51 + 0.555 + 0.35 + 0 + 0 = 0.94$$

$$\cos(\text{Doc1}, \text{Doc3}) = 0.79$$

$$\cos(\text{Doc2}, \text{Doc3}) = 0.69$$

Example 2 Doc and query

Doc: car insurance auto insurance

q: best car insurance

Term	Query		Doc
	tf - raw	df	
Auto	0	5000	1
best	1	50000	0
car	1	10000	1
insurance	1	1000	2

Product

$$= \text{Query TF-IDF} * \text{Doc TF-weight}$$

$$\text{tf-weight} = 1 + \log(\text{TF})$$

$$\text{IDF} = \log\left(\frac{\text{NT}}{\text{DF}(t)}\right)$$

Query

DAY Document
MON TUE WED THU FRI SAT SUN

Term	tf raw	tf-weight	dt	IDF	tf raw	tf-weight	Product
Auto	0	0	5000	2.3	1	1	0
best	1	1	50000	1.3	0	0	0
car	1	1	10000	2.0	1	1	2
Insurance	1	1	1000	3.0	2	1.3	3.9

$$\text{Doc length} = \sqrt{1^2 + 0^2 + 1^2 + 1 \cdot 3^2} \\ = 1.92$$

$$\text{Score} = (0+0+2+3.9) / 1.92 \\ = 3.07$$

Query TF-IDF

$$0 \cdot 2.3 = 0$$

$$1 \cdot 1.3 = 1.3$$

$$1 \cdot 2.0 = 2.0$$

$$1 \cdot 3.0 = 3.0$$

Centroid Rocchio classifier

- Text classification that involves representing each class by its centroid vector and then classify new docs based on their proximity to the class centroids
- Useful for large-scale classification with large no. of classes
- Relatively simple to implement and computationally efficient

Cons

DATE

&

Day

MON TUE WED THU FRI SAT SUN

- Assumes that each class is linearly separable from the others
- each class is represented by a single centroid vector
- May not perform well as more complex classification algs on some tasks

Large margin classifier

- Finds a hyperplane that maximizes the margin b/w different classes
- Mostly used in binary classification problems

Pros

- High Accuracy and Robustness
- Can handle noisy or overlapping data
- Applicable to wide range of classification problems
- Effective with high-dimensional data

Cons

- Computationally expensive for large datasets
- Careful selection of hyperparameters
- Not perform well on imbalanced datasets
- limited to linearly separable classes and may require complex transformations to handle non-linear decision boundaries

DATE

Day

MON TUE WED THU FRI SAT SUN

Passive Aggressive Algo

- Margin-based algo used for binary classification
- Similar to SVM but with some key differences in the way it updates its model parameters
- makes an aggressive update to the model if a data point is misclassified and passive update if correctly classified

Pros

- well-suited for online learning scenarios
- fast and efficient updates to model param
- Robust to noisy and mislabeled data

Cons

- sensitive to choice of hyperparameters
- require more iterations to converge
- limited to binary classification

Support Vector Machine

- Supervised ML Algo that works by identifying a hyperplane in a high-dimensional space that can best separate the data points into different classes
- Handle non-linear decision boundaries through the use of kernel functions

DATE



DAY

MON TUE WED THU FRI SAT SUN

Pros

- High accuracy
- Robust to overfitting
- Works well with high dimensional data
- Can handle both binary and multi-class data
- Effective in handling non-linear decision boundaries

Cons

- computationally expensive
- May not perform well with imbalanced data
- careful selection of hyperparameters
- Can be sensitive to outliers