

DWM Assignment 3

1. Total records = 14

$$H(s) = -p(\text{yes}) \log_2 p(\text{yes}) - p(\text{no}) \log_2 p(\text{no})$$

$$H(s) = 0.94$$

① outlook

$$H(\text{sunny}) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.971$$

$$H(\text{rain}) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.971$$

$$H(\text{overcast}) = -\frac{4}{4} \log_2 \left(\frac{4}{4}\right) = 0$$

$$IG = 0.94 - \left(\frac{2}{14} \times 0.971 + \frac{5}{14} \times 0.971 \right) = 0.247$$

② Temperature (i) (consider Threshold = 70)

$$H(T > 70) = -\frac{4}{9} \log_2 \left(\frac{4}{9}\right) - \frac{3}{9} \log_2 \left(\frac{3}{9}\right) = 0.991$$

$$H(T \leq 70) = -\frac{4}{5} \log_2 \left(\frac{4}{5}\right) - \frac{1}{5} \log_2 \left(\frac{1}{5}\right) = 0.721$$

$$IG = 0.94 - \left[\left(\frac{9}{14} \times 0.991 \right) + \left(\frac{5}{14} \times 0.721 \right) \right] = 0.045$$

(ii) Considering Threshold 75

$$H(T > 75) = \frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) = 0$$

$$H(T \leq 75) = \frac{7}{10} \log_2 \left(\frac{7}{10}\right) - \frac{3}{10} \log_2 \left(\frac{3}{10}\right) = 0.72$$

$$IG = 0.94 - \left(\frac{10}{14} \times 0.72 + \frac{4}{14} \times 0 \right) = 0.14$$

(iii) consider Threshold 80, we get $IG = 0.101$ IG is max for 75, so we consider it.

③ Humidity

(i) taking 80 as threshold, we get $IG = 0.101$

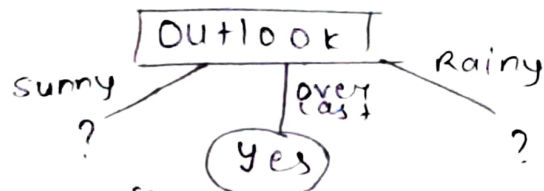
④ Wind

$$H(\text{strong}) = -\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) = 1$$

$$H(\text{weak}) = -\frac{2}{8} \log_2 \left(\frac{2}{8}\right) - \frac{6}{8} \log_2 \left(\frac{6}{8}\right) = 0.811$$

$$IG = 0.94 - \left(\frac{6}{14} \times 1 + \frac{8}{14} \times 0.811 \right) = 0.049$$

Outlook has max IG, so it is the root



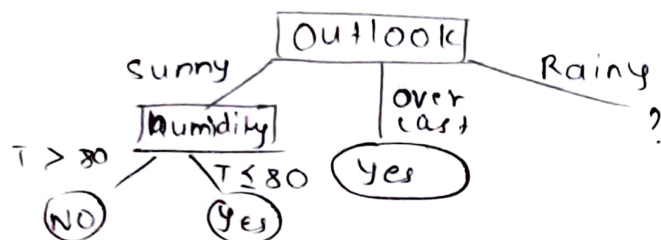
Considering outlook^{sunny}, calculate IG for all

For humidity,

$$H(T \leq 80) = -\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{2}{2}\right) = 0$$

$$H(T > 80) = -\frac{3}{3} \log_2\left(\frac{3}{3}\right) = 0$$

$\therefore IG = \max$



For Rain as outlook,

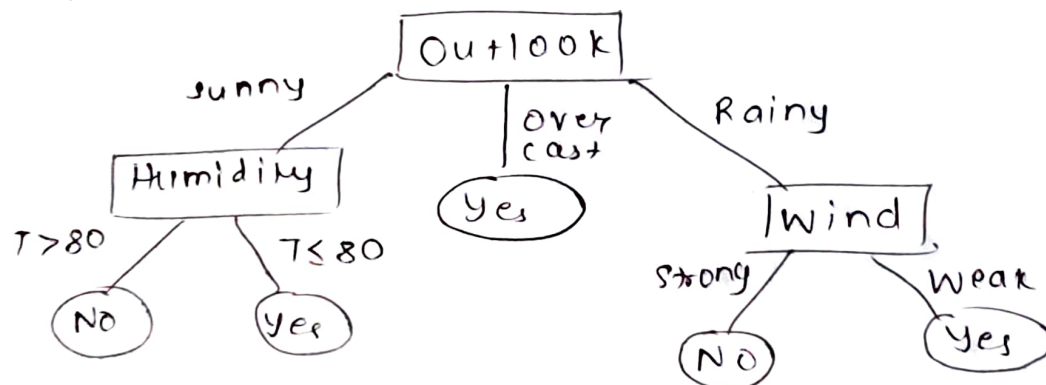
wind

$$H(\text{weak}) = 0$$

$$H(\text{strong}) = 0$$

$\therefore IG = 0.97 = \max$

\therefore Final Tree is



2. student	x	y	$x - \bar{x}$	$y - \bar{y}$	$(x - \bar{x})^2$	$(y - \bar{y})(x - \bar{x})$
1	95	85	17	8	289	136
2	85	95	7	18	49	126
3	80	70	2	-1	4	-14
4	70	65	-8	-12	64	96
5	60	70	-18	7	324	126
	<u>390</u>	<u>385</u>			<u>730</u>	<u>476</u>

$$\bar{x} = 79, \quad \bar{y} = 77, \quad y = b_0 + b_1 x$$

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{470}{730} = 0.643$$

$$b_0 = \bar{y} - b_1 \bar{x} = 77 - 0.643 \times 79 = 26.84$$

$$y = 26.84 + 0.643x$$

$$\text{when } x = 80, \quad y = 77.98$$

3. Initial value of centroids, $x_1 = A_1, x_2 = B_2, x_3 = C_1$

Iteration 0

using Euclidean Distance

D_i^0 :	A1	A2	A3	B1	B2	B3	C1	C2	(centroid)
	0	5	8.48	3.61	7.07	7.21	8.06	2.24	x_1
	3.61	4.24	5	0	3.61	4.12	7.21	1.41	x_2
	8.06	3.16	7.28	7.21	6.71	5.39	0	7.62	x_3

Object clustering

G_0 : Group 1 = A1

Grp 2 = A3, B1, B2, B3, C2

Grp 3 = A2, C1

Iteration 1, determine centroid

$$x_1 = (2, 10)$$

$$x_2 = \frac{8+5+7+6+4}{5}, \frac{4+8+5+4+9}{5} = (6, 6)$$

$$x_3 = \frac{2+1}{2}, \frac{5+2}{2} = (1.5, 3.5)$$

D1:	A1	A2	A3	B1	B2	B3	C1	C2	centroid
	0	5	8.44	3.61	7.07	7.21	8.06	2.24	x1
	5.56	4.12	2.83	2.24	1.41	2	6.40	3.61	x2
	6.52	1.58	6.52	5.7	5.70	4.52	1.58	6.04	x3

G1: Group 1 = A1, C2
 Group 2 = A3, B1, B2, B3, &
 Group 3 = A2, C1

Iteration 2: centroid

$$x_1 = \frac{2+4}{2}, \frac{10+9}{2} = (3, 9.5)$$

$$x_2 = \frac{18+5+7+6}{4}, \frac{4+8+5+4}{4} = (6.5, 5.25)$$

$$x_3 = \frac{12+1}{2}, \frac{5+2}{12} = (1.5, 3.5)$$

D2:	A1	A2	A3	B1	B2	B3	C1	C2	centroid
	1.12	2.35	7.43	2.5	6.02	6.26	7.76	1.12	x1
	6.54	4.51	1.95	3.13	0.56	1.35	6.33	7.68	x2
	6.52	1.58	6.52	5.70	5.70	4.52	1.58	6.04	x3

G2: Grp 1 = A1, B1, C2
 Grp 2 = A3, B2, B3
 Grp 3 = A2, C1

Iteration 3: centroid

$$x_1 = \frac{2+5+4}{3}, \frac{10+9+8}{3} = (3.67, 9)$$

$$x_2 = \frac{8+7+6}{3}, \frac{4+5+4}{3} = (7, 4.33)$$

$$x_3 = \frac{2+1}{2}, \frac{5+2}{2} = (1.5, 3.5)$$

D3:	A1	A2	A3	B1	B2	B3	C1	C2	centroid
	1.95	4.33	6.61	1.66	5.2	5.52	7.44	0.33	x1
	6.01	5.04	1.05	4.17	0.67	1.05	6.44	5.55	x2
	6.52	1.58	6.52	5.7	5.7	4.52	1.58	6.04	x3

G3: Grp 1 = A1, B1, C2
 Grp 2 = A3, B2, B3
 Grp 3 = A2, C1

since groups in iteration 2 and 3 are same, we stop
 .. final groups are G3

4. (i) K-Means: In k-means clustering, data objects are classified into attributes or features into k-clusters.
- Input: No of clusters (k) and points of problem
- Shape: Spherical
- Limitations: Can't cluster in geometric spaces. Gives more weight to bigger clusters than smaller ones. Overlapping clusters can't be generated.
- Outliers: The cluster centre is pushed towards the outlier, thereby distorting the actual cluster.
- (ii) K-medoids: Medoids of values in cluster are used. Attempt to minimize sum of dissimilarities between objects labelled to be in a cluster and one of objects designated as the representative of that cluster known as medoid.
- Limitation: Does not scale well, or work efficiently for large datasets
- Outliers: Partitioning among medoids is more robust than k-~~medoid~~^{means} in presence of outliers
- (iii) CLARA: Extension of k-medoids, can handle large amounts of data.
- Limitation: Doesn't guarantee perfect output for localised area.
- Uses random samples for neighbors.
- (iv) OPTICS: Ordering points to identify cluster structure. Does not explicitly segment data into clusters except it produces a visualization of reachability

distance and was distance to cluster data.

Limitations: No actual clusters, requires more memory compared to other clusters

5. candidate list = { MONKEYDACUI }

C1: Itemset	Support	Itemset	Support
M ✓	3	D	1
O ✓	3	A	1
N	2	C	2
K ✓	5	U	1
E ✓	4	I	1
Y ✓	3		

As per support = 60% or 3 transactions

L1: Itemset	Support
M	3
O	3
M	5
K	
E	4
Y	3

C2: Itemset	Support	Itemset	Support
MO	1	OE ✓	3
MK ✓	3	OY	2
ME	2	KE ✓	4
MY	2	KY ✓	3
OK ✓	3	EY	2

L2: Itemset	Support
MK	3
OK	3
OE	3
KE	4
KY	3

C3: Itemset	Support	Itemset	Support
MKO	1	OE	3
MKOE	1	OKY	2
ME	2	OEKY	2
KY	2	KEY	2

Only {OKE} satisfies support rules

Association	Confidence	
$EO \rightarrow K$	3/3	100%
$OK \rightarrow E$	3/3	100%
$KE \rightarrow O$	3/4	
$O \rightarrow KE$	3/3	100%
$E \rightarrow OK$	3/4	
$K \rightarrow EO$	3/5	

Final rules are

$$E, O \rightarrow K$$

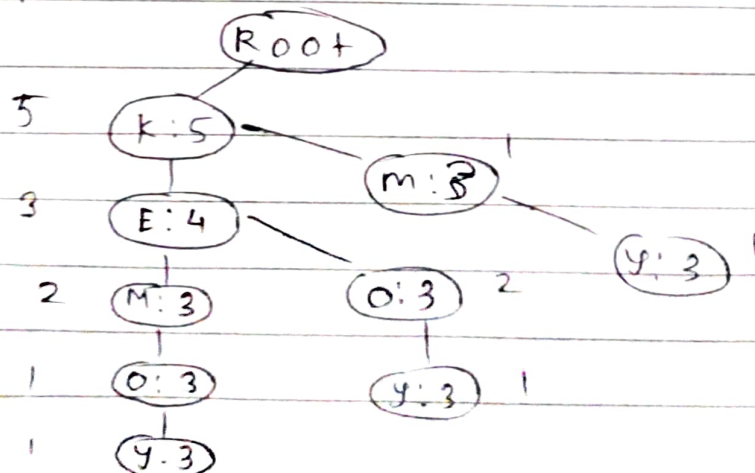
$$O, K \rightarrow E$$

$$O \rightarrow K, E$$

(iii) FP Tree : Considering C_i calculated above, we get in frequency descending order {K E m O y}

T100	MONKEY	KE m O y
T200	DONKEY	KE O y
T300	MAKE	KE m
T400	MUCKY	K m y
T500	COOKIE	KE O

FP-Tree



(iii) Vertical Data Format

A T₃
 C T₄ T₅
 D T₂
 E T₁ T₂ T₃ T₅
 I T₅
 K T₁ T₂ T₃ T₄ T₅
 M T₁ T₃ T₄
 N T₁ T₂
 O T₁ T₂ T₃
 U T₄
 Y T₁ T₂ T₄

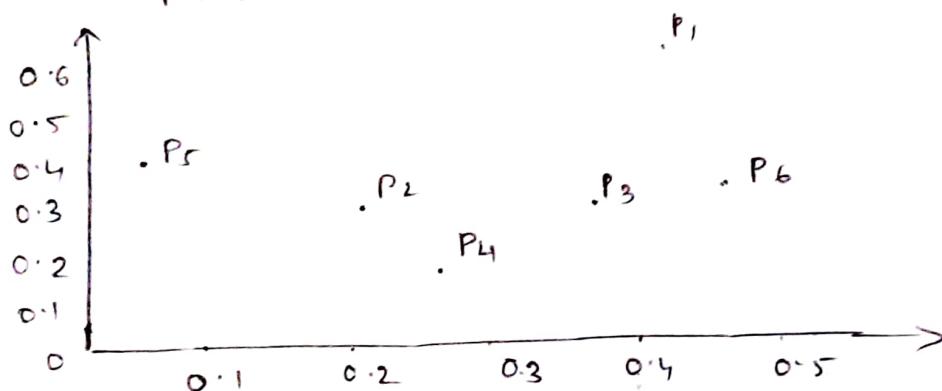
Inter
section

AE T₃
 AK T₃
 Am T₃
 CK T₄ T₅
 :
 EK T₁ T₂ T₃ T₅
 km T₁ T₃ T₄
 DK T₁ T₂ T₅
 KY T₁ T₂ T₄
 EO T₁ T₂ T₅

Inter
section

Ekm T₁ T₃
 EOK T₁ T₂ T₅
 EKY T₁ T₂

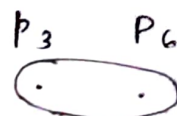
7. step 1: plot



Distance matrix

P ₁	0					
P ₂	0.24	0				
P ₃	0.22	0.15	0			
P ₄	0.37	0.2	0.15	0		
P ₅	0.34	0.14	0.28	0.29	0	
P ₆	0.23	0.25	0.19	0.22	0.39	0
	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆

(P₃, P₆) has the least distance so make them a cluster



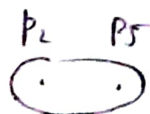
Dendrogram



Distance matrix

P ₁	0					
P ₂	0.24	0				
P ₃ P ₆	0.22	0.15	0			
P ₄	0.37	0.28	0.15	0		
P ₅	0.34	0.14	0.28	0.24	0	
	P ₁	P ₂	P ₃ P ₆	P ₄	P ₅	

(P₂, P₅) is smallest distance so make it a cluster

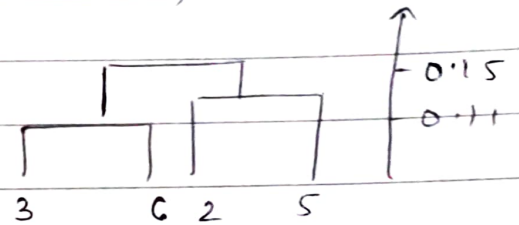


Distance matrix

P_1	0			
$P_2 P_5$	0.24	0		
$P_3 P_6$	0.22	0.15	0	
P_4	0.37	0.20	0.15	0

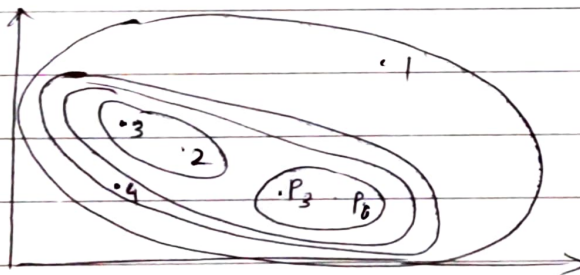
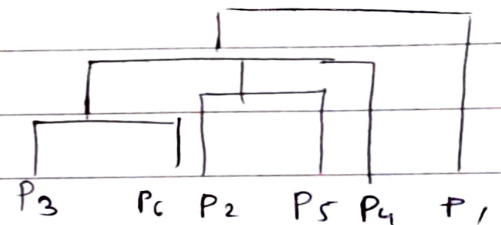
dist 0.15 is min

$\therefore P_3 P_6, P_2 P_5$ is cluster



Distance matrix

P_1	0		
$P_2 P_5 P_3 P_6$	0.22	0	
P_7	0.37	0.15	0
	P_1	$P_2 P_5 P_3 P_6$	P_4



8. We collect a set of attribute value count tables and update counts as each new example streams in. To discover evolution of classification scheme, we can maintain counts for a few classification in parallel. For instance, we can keep one classifier based on history of data, another one based on previous week of data and one based on previous days of data. For, weekly classifiers we can keep count of previous 7 days. At end of each day we discard oldest days count. For daily classifier we maintain separate counts for each hour.