Q28·14) frequent item set. & / Minimum support value = 0·2 *

| Item. | frequency | Support |
|---|---|---|
| Milk | 5 | 0·5 |
| Bread. | 4. | 0·4 |
| Eggs. | 4. | 0·4. |
| Juice | 3 | 0·3 |
| butter | 2 | 0·2 |
| Cookies. | 2 | 0·2 |
| Coffee. | 3. | 0·3. |

Since Minimum support val = 0·2
∴ all items are frequent items as all support values are greater than 0·2 ∴
and. all 1-set items are valid.

2-Item set. →

| 2-Item set | frequency → | support. |
|---|---|---|
| (Milk, Bread) | 4 | 0·4 * |
| (Milk, eggs). | 3 | 0·3 * |
| (Milk, Juice) | 1. | 0·1 |
| (bread, eggs) | 3. | 0·3. * |
| (bread, cookies). | 1 | 0·1. |
| (Juice, butter) | 1 | 0·1 |
| (cookies, butter) | 1 | 0·1 |
| (cookies, eggs). | 1 | 0·1 |
| (coffee, eggs) | 1 | 0·1 |
| (coffee, Juice) | 1 | 0·1. |

all support values less than 0·2 are not frequent items.

∴ 
| Milk, Bread | 0·4 |
| Milk, eggs | 0·3 |
| Bread, eggs | 0·3. |

3-Set items.

| | frequency | Support |
|---|---|---|
| Milk - Bread - eggs | 3 | 0·3. * |
| Milk - eggs → Bread | 0. | |
| Eggs Bread - eggs - Milk | | |

∴ 3 item set with Support > 0·2 = Milk - Bread - Eggs → 0·3

Q: 28·15 : The above question has only one frequent set of size 3.
ie Milk Eggs Bread with support 70·2.,

Milk bread eggs · = $\dfrac{Support\ (milk,\ bread,\ eggs)\ =\ 0·3}{Support\ (milk,\ bread)\ =\ 0·4}$

Milk – bread – eggs = $\dfrac{3}{4}$ = 0·75

confidence of milk eggs bread = $\dfrac{Support\ (M \cup B \cup E)}{Support\ (M \cup E)}$ = $\dfrac{0·3}{0·3}$ = $\dfrac{1}{1}$

Milk → eggs → bread = 1.

Q 28·19

Class Attribute = Repeat customer.

n ( Repeat Customer = 'yes' ) = 7
n ( Repeat Customer = 'no' ) = 3

∴ info gain = $-\left[\ 3/10\ \log_2 (3/10) + \dfrac{7}{10}\ \log_2 (7/10)\right]$

= 0·88.

→ Age attribute
① age (20–30)    freq = 5·( 4 yes    1 No).

info gain = $-\left[\ \dfrac{4}{5}\ \log (4/5) + \dfrac{1}{5}\ \log (1/5)\right]$ = 0·72

age (31–40) = freq (2) ( 1 yes    1 No)

Info gain = $-\left[\ \dfrac{1}{2}\ \log (1/2) + \dfrac{1}{2}\ \log (1/2)\right]$ = 1.

age (41–50) = freq (2)

info gain = $-\left[\ 2/2\ \log 2/2\right]$ = 0.

age (51–60) = freq (1)

info gain = 0.

$E(Age) = (0.5 \times 0.72) +. (0.2 \times 1) + 0 + 0.$

$$= 0.56.$$

$Gain(Age) = 0.88 - 0.56 = 0.32$

---

City ① NY frequency = 7 (5 yes, 2 No)

info gain = $-\left[\frac{2}{7} \log(2/7) + \frac{5}{7} \log(5/7)\right] = -\left[\frac{2}{7}(1.79) + \frac{5}{7}(0.5)\right]$

$$= 0.86$$

② LA freq = 2. (1 y, 1 No)

info gain =. $-\left[\log(1/2)\right] = 1.$

③ SF= freq = 1. ∴ Info gain = 0

$E(city) = (0.7 \times 0.86) + (0.2) + 0 = 0.8$ & $\boxed{Gain = 0.88 - 0.8 = 0.08}$

---

Gender: ① F freq (7) (2 N 5 y)

info gain = $-1\left[\frac{2}{7} \log(2/7) + \frac{5}{7} \left|\log(5/7)\right]\right] = 0.86.$

② M freq = 3.

Igain = $-1\left[\frac{2}{3} \log(2/3) + \frac{1}{3} \log(1/3)\right].$

$$= 0.92.$$

$E(Gender) = (0.7)(0.86) + (0.3)(0.92) = 0.88$

$$Gain (Gender) = 0$$

Education:

① College $freq = 6$. $(1, 5)$

Igain $= -1 [\frac{1}{6} \log(\frac{1}{6}) + \frac{5}{6} \log(5/6)] = 0.65$.

② grad $=$ $freq = 2$

$Ig = $ $\log(1) = 0$

(3) High school $freq = 2$. $\therefore$ $Ig = \log(2) = 0$

$E(education) = (0.6 \times 0.65) + 0 + 0 = 0.39$

$gain (Ed) = 0.49$

---

$gain (Edu) = 0.49$
$gain (age) = 0.32$
$gain (city) = 0.08$
$gain (gender) = 0$

Decision Tree.

**Q4** Korth & Silberschatz:

Support (Hammer) = 1/3 = 33% & Support (Nails) = 1/4 = 25%

Rule 1: ∀ Transactions T, true ⟹ buys (T, hammer)

Support = 33%   Confidence = 33%

Rule 2: ∀ Transactions T, true ⟹ buys (T, Nails)

Support = 25%   Confidence = 25%

Rule 3 ∀ transactions T, buys(T, Hammer) ⟹ buys (T, Nails)

Support = 16.5%, Confidence = $\dfrac{\text{Support(Rule 3)}}{\text{Support Rule 1}} \times 100 = \dfrac{33}{2 \times 33} \times 100$

Support = 16.5%, Confidence = 50%

Rule 4: ∀ transactions T buys(T, Nails) ⟹ buys (T, Hammer)

Support = 16.5%, Confidence = $\dfrac{S(\text{Rule 4})}{S(\text{Rule 2})} \times 100 = \dfrac{33}{200_2} \dfrac{4^2}{4} \times 100 = 66\%$

**Q5** Total Documents = 20

Total Returned Docs = 10 + 8 = 18

Precision = $\dfrac{\text{Relevant Docs}}{\text{Total Docs Returned}} = \dfrac{8}{18} = 0.44$

Recall = $\dfrac{\text{Relevant Docs}}{\text{Total Docs}} = \dfrac{8}{20} = 0.4$

# Q6. Frequency of word in following data

| Doc | computer | Doctoral | Algorithms | Watson |
|-----|----------|----------|------------|--------|
| 1 | 2 | 1 | 0 | 0 |
| 2 | 8 | 0 | 2 | 1 |
| 3 | 20 | 0 | 5 | 0 |
| 4 | 2 | 2 | 0 | 0 |
| 5 | 20 | 0 | 0 | 2 |

$$IDF = \log\left(\frac{Total\ Docs}{no\ of\ docs\ with\ term}\right) \Rightarrow IDF\ (Computer) = \log_2\left(\frac{5}{5}\right) = 0$$

$$IDF\ (Doctoral) = \log_2\left(\frac{5}{2}\right) = 1.32 \checkmark$$

$$IDF\ (Algorithms) = \log_2\left(\frac{5}{2}\right) = 1.32 \checkmark$$

$$IDF\ (Watson) = \log_2\left(\frac{5}{2}\right) = 1.32 \checkmark$$

$$TF = f_i / \left(sum\ x\ from\ 1\ to\ |v|\ f_{xj}\right)$$

$TF_{ij} \rightarrow$ normalized freq.

$f_{ij} \rightarrow$ no of occurrences $|V| \rightarrow$ features

| Doc | Computer | Doctoral | Algorithm | Watson |
|-----|----------|----------|-----------|--------|
| 1. | 2/3 = 0.67 | 1/3 = 0.33 | 0 | 0 |
| 2 | 8/11 = 0.73 | 0 | 2/11 = 0.18 | 1/11 = 0.09 |
| 3 | 20/25 = 0.8 | 0 | 5/25 = 0.2 | 0 |
| 4. | 2/4 = 0.5 | 2/4 = 0.5 | 0 | 0 |
| 5 | 20/22 = 0.91 | 0 | 0 | 2/22 = 0.09 |