

problem 1.

(a) Six items in the dataset, so, the possible rules are  $3^6 - 2^{6+1} + 1 = 602$ .

(b) 4. Since the longest transaction has 4 items.

(c)  $C_6^3 = \frac{6 \times 5 \times 4}{3 \times 2} = 20$

(d) {milk}	$\frac{5}{10}$	{milk, Diaper}	$\frac{4}{10}$
{Beer}	$\frac{4}{10}$	{milk, Bread}	$\frac{3}{10}$
{Diapers}	$\frac{7}{10}$	{milk, Butter}	$\frac{3}{10}$
{Bread}	$\frac{5}{10}$	{Diaper, Bread}	$\frac{3}{10}$
{Butter}	$\frac{5}{10}$	{Diaper, Butter}	$\frac{2}{10}$
{Cookies}	$\frac{4}{10}$	{Bread, Butter}	$\frac{5}{10}$

So {Bread, Butter} has the largest support =  $\frac{1}{2}$ .

(e) {Bread, Butter}. Since the support of {Butter}, {Bread} and {Bread, Butter} are all the same. So the confidence all equals to 1.

$$\{Bread\} \rightarrow \{Butter\} = \{Butter\} \rightarrow \{Bread\}$$



problem 2.

$$\boxed{\begin{array}{l} \{ \text{video game} \} : \{ \text{vg} \} \\ \{ \text{movie} \} : \{ \text{m} \} \end{array}}$$

we know:  $c(\{ \text{m} \} \rightarrow \{ \text{vg} \}) < \frac{\sigma(\{ \text{vg} \})}{N}$

$$\text{so: } \frac{\sigma(\{ \text{m}, \text{vg} \})}{\sigma(\{ \text{m} \})} < \frac{\sigma(\{ \text{vg} \})}{N} \quad (1)$$

$$\sigma(\{ \text{vg} \}) = \sigma(\{ \bar{\text{m}}, \text{vg} \}) + \sigma(\{ \text{m}, \text{vg} \}) \quad (2)$$

$$N = \sigma(\{ \bar{\text{m}} \}) + \sigma(\{ \text{m} \}) \quad (3)$$

$$\Rightarrow \frac{\sigma(\{ \text{m}, \text{vg} \})}{\sigma(\{ \text{m} \})} < \frac{\sigma(\{ \bar{\text{m}}, \text{vg} \}) + \sigma(\{ \text{m}, \text{vg} \})}{\sigma(\{ \bar{\text{m}} \}) + \sigma(\{ \text{m} \})} \leq 1$$

$$\Rightarrow \sigma(\{ \text{m}, \text{vg} \}) \cdot \sigma(\{ \bar{\text{m}} \}) + \cancel{\sigma(\{ \text{m}, \text{vg} \}) \cdot \sigma(\{ \text{m} \})} < \cancel{\sigma(\{ \text{m} \}) \cdot \sigma(\{ \bar{\text{m}}, \text{vg} \})} + \sigma(\{ \text{m} \}) \cdot \sigma(\{ \text{m}, \text{vg} \}) \quad (4)$$

$$\Rightarrow \frac{\sigma(\{ \text{m}, \text{vg} \})}{\sigma(\{ \text{m} \})} < \frac{\sigma(\{ \bar{\text{m}}, \text{vg} \})}{\sigma(\{ \bar{\text{m}} \})}$$

$$\text{so: } c(\{ \text{m} \} \rightarrow \{ \text{vg} \}) < c(\{ \bar{\text{m}} \} \rightarrow \{ \text{vg} \})$$

$$\text{from (4)} \Rightarrow \cancel{\sigma(\{ \text{m}, \text{vg} \}) \cdot \sigma(\{ \bar{\text{m}} \})} + \sigma(\{ \bar{\text{m}}, \text{vg} \}) \cdot \sigma(\{ \bar{\text{m}} \}) < \cancel{\sigma(\{ \bar{\text{m}} \}) \cdot \sigma(\{ \bar{\text{m}}, \text{vg} \})} + \sigma(\{ \bar{\text{m}} \}) \cdot \sigma(\{ \bar{\text{m}}, \text{vg} \})$$

$$\Rightarrow \frac{\sigma(\{ \bar{\text{m}}, \text{vg} \}) + \sigma(\{ \text{m}, \text{vg} \})}{\sigma(\{ \bar{\text{m}} \}) + \sigma(\{ \text{m} \})} < \frac{\sigma(\{ \bar{\text{m}}, \text{vg} \})}{\sigma(\{ \bar{\text{m}} \})}$$

$$\Rightarrow s(\{ \text{vg} \}) < c(\{ \bar{\text{m}} \} \rightarrow \{ \text{vg} \})$$

□



problem 3.

$\{6, 12, 18, 24, 30, 42, 48\}$

a). i.  $\{18, 45\}$ .

$$|30-18| < |30-45|.$$

So, first cluster:  $\{6, 12, 18, 24, 30\}$ .

$$\text{squared error: } \sum_{i \in C_1} (i-18)^2 = 360$$

second cluster:  $\{42, 48\}$ .

$$\text{squared error: } \sum_{i \in C_2} (i-45)^2 = 18.$$

total error: 378.

ii)  $\{15, 40\}$ .

cluster 1:  $\{6, 12, 18, 24\}$ , error = 180

cluster 2:  $\{30, 42, 48\}$ , error = 168

total error: 348

b) yes, both are stable,

Since the average for each cluster in each set equals the initial centroid value set

$$\text{e.g.: } \frac{6+12+18+24+30}{5} = 18; \quad \frac{42+48}{2} = 45$$

c). the two clusters are  $\{6, 12, 18, 24, 30\}$  &  $\{42, 48\}$

d) single link produces more natural clustering.  
since clusters generated by single link has uniform distance between nodes.

e) Contiguous and also density.

single link:  $\{6, 12, 18, 24, 30\}$  &  $\{42, 48\}$

k-means:  $\{6, 12, 18, 24\}$  &  $\{30, 42, 48\}$



f). the well-known drawback of k-means is its bad performance at dividing the clusters with big-different-sizes.

the k-mean algorithm tends to break the larger cluster as we can see from this example.