
Data Mining Exercises

1 Krimp

In the code tables below, only the length of the codes of the itemsets are given (in bits), not the codes themselves (since we are not interested in the codes).

| Tid | Transaction |
|-----|-------------|
| 1 | ABCD |
| 2 | BCD |
| 3 | ACD |

| Itemset | Code Length |
|---------|-------------|
| ABC | 1.5 |
| BC | 2 |
| CD | 2.5 |
| C | 3.5 |
| A | 2.5 |
| D | 3 |

| Tid | Transaction |
|-----|-------------|
| 1 | ABC |
| 2 | ABC |
| 3 | ABC |
| 4 | BC |
| 5 | BCD |
| 6 | CD |

~~Question 1~~ : Compute $\text{Len}(DB_1 \mid CT_1)$. For each transaction, give the itemsets used to compress it.

We assume the input of Krimp is DB_2 and the set of candidate itemsets is $F = \{ABC, BC, CD\}$.

~~Question 2~~ : In which order are the candidate itemsets considered by Krimp?

~~Question 3~~ : Give the Standard Code Table ST for DB_2 and compute its length $\text{Len}(ST)$.

~~Question 4~~ : Compute $\text{TotLen}(DB_2, ST) = \text{Len}(DB_2 \mid ST) + \text{Len}(ST)$

~~Question 5~~ : Suppose that a code table CT_2 contains itemsets A, B, C, D, BC, ABC. This code table is used to compress DB_2 . In which order do the itemsets appear in the code table? Compute the code lengths of each of these itemsets in CT_2 (if an itemset has a usage of 0, its length is 0). You can use the logarithm table at the end of this subject.

~~Question 6~~ : Compute $\text{TotLen}(DB_2, CT_2) = \text{Len}(DB_2 \mid CT_2) + \text{Len}(CT_2)$

2 Sequences

| Sid | Sequence |
|-----|-----------------|
| 1 | abc(ab)bc |
| 2 | (abc)baa(abc) |
| 3 | abcc(ab)c |
| 4 | (abc)(abc)(abc) |

Question 7 : Of which sequences of the database are sequences $S1 = (ab)(bc)$ and $S2 = b(ab)c$ subsequences? You have to give the correspondence between the items of $S1$ and $S2$ and the items of the sequences of the database. What are the support of $S1$ and $S2$ in this database?

Question 8 : Given the following set of frequent sequences of size 3, what are the candidates sequences of size 4? How do you compute them? $F_3 = \{a(ab), (ab)c, aac, abc, bcb, (ab)b, (abc)\}$

3 Base 2 Logarithm Table

Recall that $\log_2(ab) = \log_2(a) + \log_2(b)$ and thus $\log_2(a/b) = \log_2(a) - \log_2(b)$.

| | | | | | | | | | | | | | | | | |
|-------------|---|---|-----|---|-----|-----|-----|---|-----|-----|-----|-----|-----|-----|-----|----|
| x | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| $\log_2(x)$ | 0 | 1 | 1.6 | 2 | 2.3 | 2.6 | 2.8 | 3 | 3.2 | 3.3 | 3.5 | 3.6 | 3.7 | 3.8 | 3.9 | 4 |

Question 1

| | | | |
|---|--|---------------------------------|-----|
| | | $A \underline{Bc}$ | 1.5 |
| 1 | | $\underline{ABC} \underline{D}$ | 2 |
| 2 | | $\underline{BC} \underline{CD}$ | 2.5 |
| 3 | | $\underline{A} \underline{CD}$ | 3.5 |
| | | D | 2.5 |
| | | | 3 |

$$\begin{aligned} \text{len}(DB_1 | CT_1) &= 1 \times \text{code}_{CT}(ABC) + 2 \times \text{code}_{CT}(D) \\ &\quad + 1 \times \text{code}_{CT}(BC) + 1 \times \text{code}_{CT}(A) + \\ &\quad 1 \times \text{code}_{CT}(CD) \\ &= 1.5 + 2 \times 3 + 2 + 2.5 + 2.5 \\ &= 1.5 + 6 + 2.5 \\ &= 14.5 \end{aligned}$$

2. $F = \left\{ \begin{array}{l} ABC, BC, CD \\ 3 \quad 5 \quad 2 \end{array} \right\}$

The order is $\left\{ \begin{array}{l} BC, ABC, CD \\ 5 \quad 3 \quad 2 \end{array} \right\}$

| | |
|---|-------------------|
| 1 | \underline{ABC} |
| 2 | \underline{ABC} |
| 3 | \underline{ABC} |
| 4 | \underline{BC} |
| 5 | \underline{BCD} |
| 6 | \underline{CD} |

3.

| X | Usage | Code length |
|-------|-------|-------------|
| A | 3 | |
| B | 5 | |
| C | 6 | |
| D | 2 | |
| Total | 16 | |

| X | Usage | Code length |
|-------|-------|------------------------------|
| C | 6 | $-\log_2 \frac{6}{16} = 1.4$ |
| B | 5 | $-\log_2 \frac{5}{16} = 1.7$ |
| A | 3 | $-\log_2 \frac{3}{16} = 2.4$ |
| D | 2 | $-\log_2 \frac{2}{16} = 3$ |
| Total | 16 | |

$$\text{len}(ST) = \sum_{c \in ST} \text{len}(c, \text{code}_{ST}(c))$$

$$= \sum_{c \in ST} \text{len}(c) + \text{len}(\text{code}_{ST}(c))$$

$$= 2 \times (1.4 + 1.7 + 2.4 + 3) = 17$$

$$\text{len}(DB_2 | ST) = 6 \times 1.4 + 5 \times 1.7 + 3 \times 2.4 + 2 \times 3$$

$$= 30.1$$

$$\text{Total len}(DB_2, ST) = 30.1 + 17 = 47.1$$

$CT_2 =$

| X | Usage | Code length |
|-------|-------|-----------------------------|
| ABC | 3 | $-\log_2 \frac{3}{8} = 1.4$ |
| BC | 2 | $-\log_2 \frac{2}{8} = 2$ |
| C | 1 | $-\log_2 \frac{1}{8} = 3$ |
| B | 0 | 0 |
| A | 0 | 0 |
| D | 2 | $-\log_2 \frac{2}{8} = 2$ |
| Total | | 8.4 |

| Tid | Transaction |
|-----|-------------|
| 1 | ABC |
| 2 | ABC |
| 3 | ABC |
| 4 | BC |
| 5 | BCD |
| 6 | CD |

$$\begin{aligned} \text{len}(ABC) &= \text{len}(\text{code}_{ST}(A)) \\ &\quad + \text{len}(\text{code}_{ST}(B)) + \text{len}(\text{code}_{ST}(C)) \\ &= 2.4 + 1.7 + 1.4 = 5.5 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{len}(BC) &= \text{len}(\text{code}_{ST}(B)) + \text{len}(\text{code}_{ST}(C)) \\ &= 1.7 + 1.4 = 3.1 \text{ bits} \end{aligned}$$

$$\text{len}(C) = 1.4 \text{ bits}$$

$$\text{len}(D) = 3 \text{ bits}$$

$$\begin{aligned} \text{len}(CT_2) &= 8.4 + 3 + 1.4 + 5.5 + 3 + 1 \\ &= 8.4 + 6.9 + 3.1 + 3 \\ &= 8.4 + 10 + 3 = 21.4 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{len}(DB_2 | CT_2) &= 3 \times 1.4 + 2 \times 2 + 3 \times 1 + 2 \times 2 \\ &= 4.2 + 4 + 3 + 4 \\ &= 15.2 \text{ bits} \end{aligned}$$

$$\text{Total len}(DB_2 | CT_2) = 21.4 + 15.2$$

$$= 36.6$$

Database : $[\underline{ABC}, \underline{ABC}, \underline{ABC}, \underline{BC}, \underline{BCD}, \underline{CD}]$

F : $\begin{bmatrix} BC & ABC & C D \\ 5 & 3 & 2 \end{bmatrix}$

Item: $[A, B, C, D]$

Candidate is BC

Code Table

| Item | usage | Code Length | length |
|------|-------|-------------|--------|
| BC | 5 | 1.1 | 3.1 |
| C | 1 | 3.5 | 1.4 |
| B | 0 | 0 | 0 |
| A | 3 | 1.9 | 2.4 |
| D | 2 | 2.5 | 3.0 |

$$\text{Total} = (1.1 + 3.5 + 1.9 + 2.5 + 3.1 + 1.4 + 2.4 + 3.0) \\ = 18.900$$

Compressed Database

| | | |
|---|--------------------------------|-----|
| 1 | $\underline{A} \underline{BC}$ | 3.0 |
| 2 | $\underline{A} \underline{BC}$ | 3.0 |
| 3 | $\underline{A} \underline{BC}$ | 3.0 |
| 4 | \underline{BC} | 1.1 |
| 5 | $\underline{BC} \underline{D}$ | 3.6 |
| 6 | \underline{CD} | 6.0 |

Total database

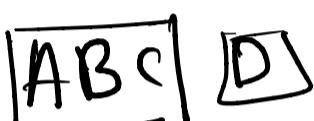
$$\text{length} = (18.900 + 19.7) = 38.6$$

*** candidate OK

Candidate is ABC

| TID | | ItemSet | Code Length |
|-----|------|---------|-------------|
| 1 | ABCD | ABC | 1.5 |
| 2 | BCD | BC | 2 |
| 3 | ACD | CD | 2.5 |

| CT ₁ |
|-----------------|
| C |
| A |
| D |
| B |

1 

$$\text{len}(DB_1)_{CT_1}$$

2 

$$= 1.5 + 2 \times 3 + 2 + 2.5 + 2.5$$

③ 

$$= 8.5 + 6 = 14.5$$

②

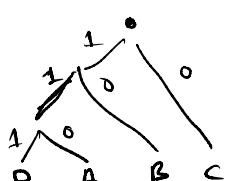
| Tid | Transaction |
|-----|-------------|
| 1 | ABC |
| 2 | ABC |
| 3 | ABC |
| 4 | BC |
| 5 | BCD |
| 6 | CD |

$$F = \{ABC, BC, CD\}$$

By Krimph's algo, $F' = \{ABC, BC, CD\}$

| X | Code(x) | occur |
|-------|---------|-------|
| A | 110 | 3 |
| B | 10 | 5 |
| C | 0 | 6 |
| D | 111 | 2 |
| Total | | 16 |

$$\begin{aligned} \text{len}(ST) &= \sum_{x \in ST} \text{Code}_{ST}(x) \\ &= -\left(\log \frac{3}{16} + \log \frac{2}{16} + \log \frac{6}{16} + \log \frac{5}{16}\right) \end{aligned}$$



$$= -\left(\log \frac{1 \cdot 5 \cdot 6 \cdot 2}{16 \times 16 \times 16 \times 16} \right) = -\log \left(\frac{15}{16^3 \times 2} \right) = 2.74$$

$$\ln(DB_2/ST) = - \left(3 \cdot \log \frac{3}{16} + 5 \cdot \log \frac{5}{16} + 6 \cdot \log \frac{6}{16} + 2 \cdot \log \frac{2}{16} \right)$$

$$= 9.07$$

$$\text{total}(DB_2/ST) = 9.07 + 2.74$$

$$= 11.81$$

A, B, C, D, BC, ABC
 3 5 6 2, 5 3

| | Itemset | Support | $\ln(\text{Code}_{ST}(x))$ |
|-------------|---------|---------|----------------------------|
| C_{T_2} = | ABC | 3 | $-\ln \frac{3}{8}$ |
| | BC | 2 | $-\ln \frac{2}{8}$ |
| | C | 1 | $-\ln \frac{1}{8}$ |
| | BA | 0 | 0 |
| | A | 0 | 0 |
| | D | 2 | $-\ln \frac{2}{8}$ |
| | Total | 8 | |

$$\ln(C_{T_2}) = -\log \frac{3 \cdot 2 \cdot 1 \cdot 2}{8 \cdot 8 \cdot 8 \cdot 8} + \sum \ln(\text{Code}_{ST}(x))$$

$$= 2.55 + \sum_{x \in C_{T_2}} \ln(\text{Code}_{ST}(x))$$

$$= 2.55 - \left(\lg \frac{3}{16} + \lg \frac{5}{16} + \ln \frac{6}{16} \right) - \left(\lg \left(\frac{5}{16} \right) + \lg \frac{6}{16} \right)$$

$$+ \left(\ln \frac{6}{16} + \ln \frac{2}{16} \right)$$

$$= 2.52 + 1.66 + 0.93 + 1.33$$

$$= 6.47$$

| Tid | Transaction |
|-----|-------------|
| 1 | ABC |
| 2 | ABC |
| 3 | ABC |
| 4 | BC |
| 5 | BCD |
| 6 | CD |

$$\text{len}(DB_2 / CT_2) = - \left(3 \times \lg \frac{3}{8} + 2 \times \lg \frac{2}{8} + \ln \frac{1}{8} + 2 \times \ln \frac{2}{8} \right)$$

$$= 4.59$$

$$\text{Total len}(DB_2 / CT_2) = 6.47 + 4.59$$

$$= 11.06$$