**Data Mining (CSE542)**

**Homework 05**

ID: ____     Name: ___조원석__     Date:___2023.04.17__

**Task-1**

Consider the database shown in Table 10.2. Answer the following questions:
**(a)** Let *minsup* = 4. Find all frequent sequences.
**(b)** Given that the alphabet is $\Sigma = \{A, C, G, T\}$. How many possible sequences of length $k$ can there be?

Table 10.2. Sequence database

| Id | Sequence |
|---|---|
| $s_1$ | *AATACAAGAAC* |
| $s_2$ | *GTATGGTGAT* |
| $s_3$ | *AACATGGCCAA* |
| $s_4$ | *AAGCGTGGTCAA* |

**(a)** Let minsup = 4. Find all frequent sequences.
A−4, G−4, T−4
AA −4, AG −4, AT −4, GA −4, TA −4, TG −4
AAT −4, AGA −4, ATA −4, ATG −4, GAA −4, TAA −4, TGA−4,
AATA − 4, ATGA − 4

**(b)** Given that the alphabet is = {A,C,G,T}. How many possible sequences of length k can there be?
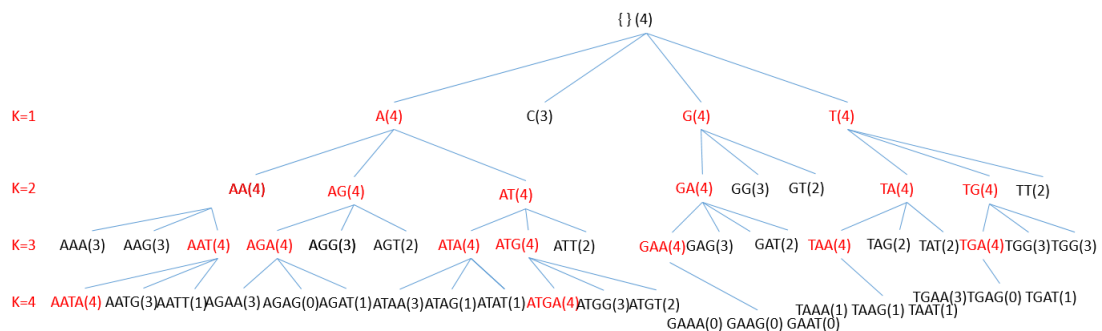$4^k$
$k = 1, 4^1 = 4$
$k = 2, 4^2 = 16$
$k = 3, 4^3 = 64$
$k = 4, 4^4 = 256$

**(c)** Show the steps of the PrefixSpan algorithm

**Task-2**

(a) Let minsup = 4. Find all frequent sequences

(b) Show the steps of the PrefixSpan algorithm

**Table 10.3.** Sequence database

| Id | Sequence |
|---|---|
| $s_1$ | *ACGTCACG* |
| $s_2$ | *TCGA* |
| $s_3$ | *GACTGCA* |
| $s_4$ | *CAGTC* |
| $s_5$ | *AGCT* |
| $s_6$ | *TGCAGCTC* |
| $s_7$ | *AGTCAG* |

**(a)** Let minsup = 4. Find all frequent sequences

A−7, C−7, G−7, T−7

AC −6, AG −6, AT −6, CA −6, CC −4, CG −6, CT −5,

GA −5, GC −6, GG −4, GT −6, TA −5, TC −6, TG −5,

ACT −4, AGC−6, AGT −5, ATC−5, CAG −4, CGC −4, CTC −4,

GAG −4, GCA−4, GCG −4, GTC −5, TCA −4, TCG −4,

AGTC −4

**(b)** Show the steps of the PrefixSpan algorithm

{ } (7)

K=1    A(7)    C(7)    G(7)    T(7)

K=2    AA(3)  AC(6)  AG(6)   AT(6)    CA(6)  CC(4)  CG(6)  CT(5)    GA(5)  GC(6)  GG(4)  GT(6)    TA(5)  TC(6)  TG(5)  TT(1)

K=3    ACA(3) ACG(3) AGA(3) AGG(2) ATA(3) ATG(3) CAA(0) CAG(4) CCA(2) CCG(0) CGA(3) CGG(1) CTA(2) CTG(2) GAA(1) GAG(4) GCA(4) GCG(4) GGA(1) GGG(1) GTA(3) GTG(3) TAA(0) TAG(3) TCA(4) TCG(4) TGA(3) TGG(1)
       ACC(3) ACT(4) AGC(6) AGT(5) ATC(5) ATT(0) CAC(2) CAT(2) CCC(2) CCT(1) CGC(4) CGT(3) CTC(4) CTT(0) GAC(3) GAT(2) GCC(3) GCT(3) GGC(1) GGT(1) GTC(5) GTT(0) TAC(2) TAT(1) TCC(2) TCT(1) TGC(2) TGT(1)

K=4    ACTA(0) AGCA(0) AGTA(2) ATCA(1)    CAGA(0)    CGCA(1) CTCA(2)    GAGA(1) GCAA(0) GCCA(1)    GTCA(3)    TCAA(0) TCGA(1)
       ACTC(3) AGCC(3) AGTC(4) ATCC(0)    CAGC(1)    CGCC(2) CTCC(1)    GAGC(2) GCAC(2) GCCC(1)    GTCC(1)    TCAC(2) TCGC(1)
       ACTG(1) AGCG(1) AGTG(2) ATCG(2)    CAGG(0)    CGCG(0) CTCG(0)    GAGG(0) GCAG(3) GCCG(1)    GTCG(2)    TCAG(3) TCGG(0)
       ACTT(0) AGCT(0) AGTT(0) ATCT(0)    CAGT(2)    CGCT(1) CTCT(0)    GAGT(1) GCAT(1) GCCT(1)    GTCT(0)    TCAT(1) TCGT(1)