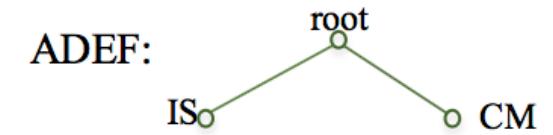
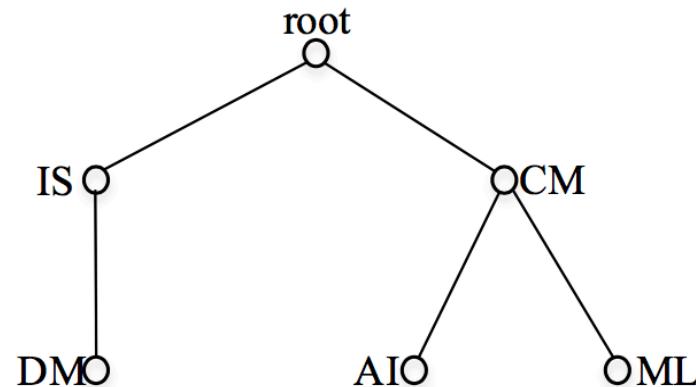
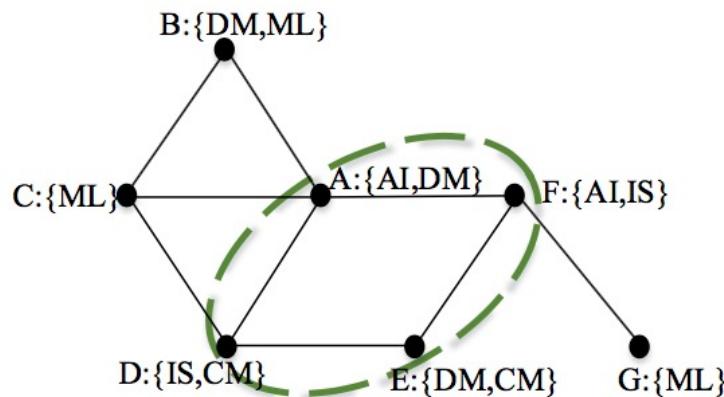


Hardness and compressed index

6.26

Problem definition

- Given a profiled graph G , a profile tree T , and one query node $q \in G$, find communities, each of which is a subgraph G_q satisfying the following properties:
 - Connectivity: G_q is connected and $G_q \subseteq G$;
 - Structure cohesiveness: k -core;
 - Semantic cohesiveness: all the vertices in G_q share the maximal common induced rooted subtree of q ;



Transformation from tree to sequence

- Semantic cohesiveness: all the vertices in G_q share the maximal common induced rooted subtree of q ;
- Mining maximal common subtree \longleftrightarrow mining maximal common subsequence \longleftrightarrow mining maximal frequent subsequence:
 - (1) $\text{min_sup}=k+1$
 - (2) subsequence of query vertex's.

Hardness of mining all frequent subsequences

- “If a counting problem is #P-complete[3] or #P-hard, then its associated problem of enumerating all solution must be NP-hard.”[1,2]
- Our target is to prove that counting the number of all maximal common subsequences is #P-hard.

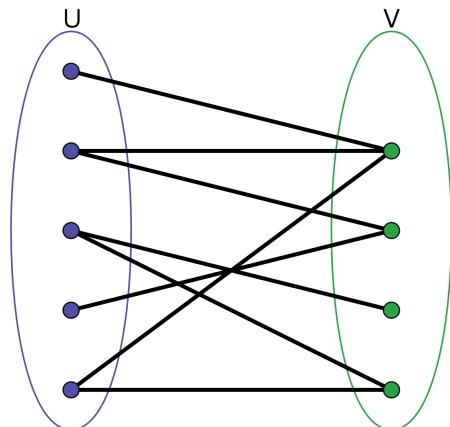
[1].Michael R G, David S J. Computers and intractability: a guide to the theory of NP-completeness[J]. WH Free. Co., San Fr, 1979: 90-91.

[2].Papadimitriou C H. Computational complexity[M]. John Wiley and Sons Ltd., 2003.

[3].Valiant L G. The complexity of computing the permanent[J]. Theoretical computer science, 1979, 8(2): 189-201.

Bipartite graph

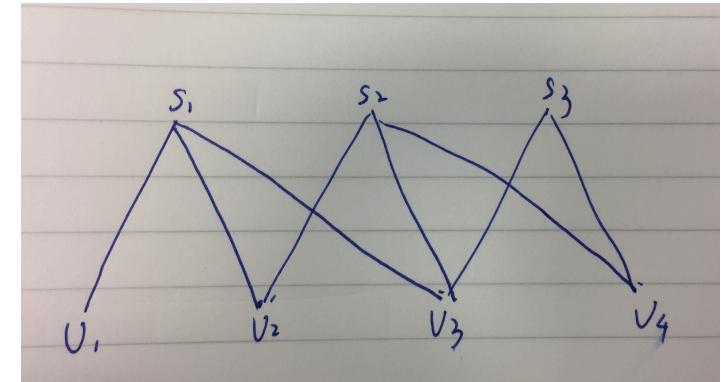
- $G = \{U, V, E\}$, U, V are two disjoint vertex sets, and E represent edges between two vertices in U and V .
- Bipartite clique: $G = \{U, V, E\}$. If there is an edge between every pair of vertices in U and V . G is bipartite clique.
- Maximal bipartite clique.
- “The problem of counting the number of maximal Bipartite cliques in a given bipartite graph is #P-complete”.[1]



[1]Provan J S, Ball M O. The complexity of counting cuts and of computing the probability that a graph is connected[J]. SIAM Journal on Computing, 1983, 12(4): 777-788.

Cont.

- Let $S = \{s_1, s_2, \dots, s_n\}$ represents all items of P-trees. $U = \{u_1, u_2, \dots, u_n\}$ represents all users in G .
- mining maximal subsequences \longleftrightarrow mining all maximal bipartite cliques.
- However, “common”(frequency) has not been included.



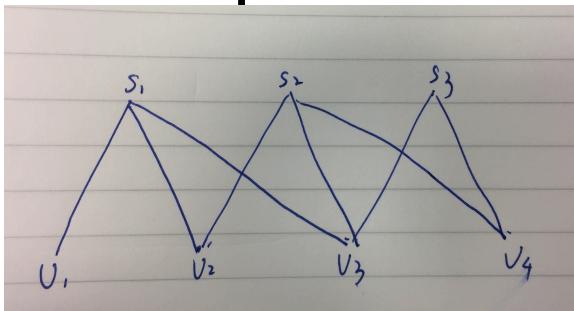
Cont.

- Let $F_i(D)$ denote all maximal frequent subsequences of D whose support is i .
- Let $M_i(D)$ denote all maximal frequent subsequences of D whose support is no less than i .
- mining maximal subsequences \longleftrightarrow mining all maximal bipartite cliques.
- all maximal bipartite cliques are $\sum_{i=0}^{\max} F_i(D)$.
- Corollary: The problem of counting the number of $\sum_{i=0}^{\max} |F_i(D)|$ is #P-complete.

Cont.

- Lemma: It is a #P-hard to count the number of all maximal common subsequences.
- Proposition: $F_k(D) \subseteq M_k(D) \subseteq \sum_{i=k}^{\max} F_i(D)$.
- Our target is to create a new database N_D and $|M_k(N_D)|$ has equation with $\sum_{i=k}^{\max} F_i(D)$. Then we can reduce *Corollary* to lemma and prove counting the number of all maximal common subsequences is at least hard as *Corollary*.

Matrix representation



	s_1	s_2	s_3
v_1	1	1	0
v_2	1	1	0
v_3	0	1	1
v_4	0	1	1

- $N_D \text{ upper} = D_q$ induced by query Vertex+ X ;
- $N_D \text{ lower} = S_Q + X - \{X_i\}$;

N_D	s_{b_1}	s_{b_2}	s_{b_3}	\dots	s_{b_l}	x_1	$x_2 \dots x_n$
U_{a_1}	1				1	1	
U_{a_2}	0	1	-	-	1	1	1
U_{a_3}					1	1	1
\vdots							
U_{a_m}							
Z_1	1	1	1	1	1	0	1
Z_2						1	0
\vdots						1	0
Z_m	1	1	1	1	1		

- Lemma: If $I \subseteq S_q$ is a maximal $(a+k)$ -support sequence in D_q and $X_k \subseteq X$. Then $I \cup X_k$ is a maximal $(a+m)$ -support sequence in N_D
- Theorem: sequence U is a maximal $(a+m)$ common subsequence in N_D iff U 's support is $(a+m)$ in N_D .

N_D	U_{a1}	S_{b1}	S_{b2}	S_{b3}	\dots	S_{bl}	X_1	X_2	\dots	X_n
	U_{a2}									
	U_{a3}									
	\vdots									
	U_{am}									
	Z_1	1	1	1	1	1	0	1	1	
	Z_2						1	0	1	
	\vdots						{	0	-	
	Z_m	1	1	1	1	1				

proof

- Theorem: sequence U is a maximal $(a+m)$ common subsequence in N_D iff U 's support is $(a+m)$ in N_D .
 - \leftarrow from definition.

3 | 72

If $1 \leq s_k$ is a maximal $(\ell+k)$ -support sequence in D_ℓ ,
 and $X_k \subseteq X$, $|X_k| = k$. We claim $\{X_k\}$ is a maximal
 $(\ell+m)$ -support sequence.

$$\sup([]) \text{ in } D_2 = \partial + k$$

$$\begin{aligned} \text{Sup}(2) \text{ in } ND &= \text{Sup}(2) \text{ in } ND^{\text{upper}} + \text{Sup}(2) \text{ in } ND^{\text{lower}} \\ &= j+k+m-k \\ &= m+j \end{aligned}$$

i. for any k ($k \leq m-1$) $(U \cup X(k))$ is a maximal
 $(\geq m)$ support sequence.

充要性: $U = \bigcup_{k=1}^{\infty} X_k$, $\exists q$ s.t. $X_k \subseteq X$. ($X_{k+1} = k$)

$\therefore U$ is maximal $(2+m)$ supporting in ND.

$\therefore \text{freq}(v) \text{ in } NP \geq 2^{rm}$.

$$\begin{aligned} \text{In } ND_{\text{lower}}: \quad fre(V) &= fre(X_k) = m-k \\ \therefore fre(V) \text{ in } ND_{\text{upper}} &= fre(V) - (m-k) \\ &\geq d+m - (m-k) \\ &= d+k. \end{aligned}$$

~~2p~~ → ND upper

- support segment. ($k' \geq k$). (2)

由引理 $\|U\|_{\infty} = \max_{k=1}^m \|U_k\|$ 是一个最大值。

(2) m) - support
它會慢慢退

$N.D.$	S_{b_1}	S_{b_2}	$S_{b_3} \dots$	S_{b_k}	X_1	$X_2 \dots X_n$
U_{a_1}	1	1	0	1	1	1
U_{a_2}	0	1	-	-	1	1
U_{a_3}	1	-	-	-	1	1
U_{a_m}	-	-	-	-	-	-
Z_1	1	1	1	1	0	1
Z_2	1	-	-	-	1	1
Z_m	1	1	1	1	-	-

Conclusion

- Theorem: sequence U is a maximal $(\sigma+m)$ common subsequence in N_D iff U 's support is $(\sigma+m)$ in N_D .

$$|M_{a+m}(N_D)| = |F_{a+m}(N_D)|$$

Lemma: If $I \subseteq S_q$ is a maximal $(a+k)$ -support sequence in D_q and $X_k \subseteq X$. Then $I \cup X_k$ is a maximal $(a+m)$ -support sequence in N_D .

$$|F_{a+m}(N_D)| = \sum_{i=0}^{m-a} |F_{a+i}(D)| \cdot C_m^i.$$

Then: $|M_{a+m}(N_D)| = \sum_{i=0}^{m-a} |F_{a+i}(D)| \cdot C_m^i.$

Reduction

- $|M_{a+m}(N_D)| = \sum_{i=0}^{m-a} |F_{a+i}(D)| \cdot C_m^i.$

$$\begin{bmatrix} F_1 \\ F_2 \\ F_3 \\ F_4 \\ \vdots \\ \vdots \\ F_m \end{bmatrix} = \left(\begin{array}{cccccc} C_m^0 & C_m^1 & \dots & C_m^{m-1} \\ 0 & C_m^0 & \dots & C_m^{m-2} \\ 0 & 0 & C_m^0 & \dots & C_m^{m-3} \\ \vdots & & & & \vdots \\ & & & & 1 \\ & & & & & 1 \end{array} \right)^{-1} \begin{bmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ \vdots \\ \vdots \\ u_m \end{bmatrix}$$

- If there exists a polynomial-time algorithm to count the number of maximal common subsequences, then $|M_{a+m}(N_D)|$ can be computed in polynomial time, then (F_1, F_2, \dots, F_m) can be computed in polynomial time.
- Corollary: The problem of counting the number of $\sum_{i=0}^{\max} |F_i(D)|$ is #P-complete.
- Then Lemma: It is a #P-complete to count the number of all maximal common subsequences.

compressed index-mathematic preliminary

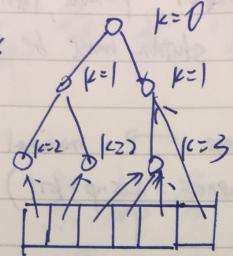
- The only divisors of a prime integer p ($p > 1$) are 1 and p .
- Every positive integer n is either 1 or can be expressed as a product of several prime integers, and this factorization is unique with the order of prime integers. The *standard form* of n factorization of n :
 $n = p_1^{m_1} p_2^{m_2} \cdots p_n^{m_n}$, p_i is a distinct prime integer, m_i is called the *multiplicity* of p_i .
- Give two integers a and b , the great common divisor of a, b is $\gcd(a, b)$. E.g., $a = 2^3 \cdot 3^2 \cdot 7 = 504$, $b = 2^2 \cdot 3^1 \cdot 7 \cdot 11 = 924$, $\gcd(a, b) = 2^2 \cdot 3^1 \cdot 7 = 84$.
- if we set m_i as 1. Then $n = p_1 p_2 \cdots p_n$. $\gcd(a, b) = \prod_{i=1}^m p_{x_i}$.

compressed index-example

- $G = \{2,3,5,7\}$, $a = 2 \cdot 3 = 6$, $b = 5 \cdot 7 = 35$, $\gcd(a,b)=1$.
- $S_a = 1100$, $S_b = 0011$, $S_a \cap S_b = 0000 = 1$.
- Set 4 bits as a block. $S_a = \{10, \dots\}$.
Then $S_a \cap S_b = \{\gcd(S_{a_i}, S_{b_i}) \mid i \in \text{number of blocks}\}$.
- If $G=\{2,3,5,7\}$, then $\gcd(a,b)$ has $2^4 \cdot 2^4 \cdot 0.5 = 128$ types which can be pre-computed and stored in a table.

Index 3

k-core index



Map[LD, Node]

p-tree index.

a $\{10:35\}, \{4:6\}, \{8:15\}$
 $G = \{2, 3, 5, 7\}$

b $\{10:35\}, \{4:2\}, \{8:15\} \cup \{10:7\}, \{4:6\}, \{8:15\}$

0011 = 35

1100 = 6

0110 = 15

$3 \times 4 (0)$

$\{1:35\}, \{4:6\}, \{8:15\} : 0011|00001000|1100|0000.000|0110$
 0 1 2 3 4 5=7 8

代表 "ai" 之皮 $\Rightarrow 3, 4, 17, 18, 33, 34$ 共有。

space cost: $(n + \frac{1}{4} \cdot n) = (\frac{5}{4}) \cdot n$.

i : the size of p-trees of n vertices