

EE 3980 Algorithms

Homework 7. Grouping Friends

105061110 周柏宇

2020/4/26

1. Introduction

In this homework, we try to find a group of people in which they can communicate directly or indirectly through other members in the same group. We first argue that this problem is essentially the same as finding the strongly connected components in a directed graph. Next, we analyze the algorithms used to solve the problem. To validate our derivation, we run the program with a different number of people and communication records to see the trend of CPU time growth.

2. Analysis & Implementation

2.1 Grouping Friends and Strongly Connected Components

In this homework, we are given the data which contains a list of names and communication records (name1 \rightarrow name2). The objective is to form the groups such that the member of the group satisfies the following rules:

1. They have sent and received messages directly between them.
2. They have sent and received messages through one or more friends

between them.

If we think of the people as the vertices and the communication records as the directed edges. Then the rules are equivalent to that the members of the group are mutually reachable. Such set of vertices are called strongly connected components.

2.2 Strongly Connected Components

Here is a high-level description of the algorithm which can be used to find the strongly connected components

```
1. // to find the strongly connected components of the
2. // graph  $G = (V, E)$ 
3. // input:  $G$ 
4. // output: strongly connected components
5. Algorithm SCC( $G$ )
6. {
7.     Construct the transpose graph  $G^T$ ;
8.     DFS_Call( $G$ ); // perform DFS to get array  $f[1 : |V|]$ 
9.     Sort  $V$  of  $G^T$  in order of
10.         decreasing value of  $f[v]$ ,  $v \in V$ 
11.     DFS_Call( $G^T$ );
12.     Each tree of the resulting DFS forest is
13.         a strongly connected components
14. }
```

As we can see, the algorithm consists of 4 parts. First, we need to construct the transposed version of the graph G^T . Second, we traverse the graph G to get the finish order. Next, we use the finish order to sort the vertices of G^T . Using the

property that members in the strongly connected components will be the same in G^T , we traverse G^T to get a forest of DFS trees. Each resulting tree is a strongly connected component.

2.3 Depth First Search

```
1. // initialization and recursive DFS function call
2. // input: graph G
3. // output: f[|V|]: finish order
4. //          SCCs: resulting DFS forest
5. Algorithm DFS_Call(G)
6. {
7.     for v := 1 to |V| do {
8.         visited[v] := 0; // initialize to not visited
9.         f[v] := 0; // reset finish order
10.    }
11.    time := 0; // global variable to track time
12.    for v := 1 to |V| do { // to handle forest case
13.        if (visited[v] = 0) then {
14.            SCCs[ptr] := -1; // separate trees
15.            ptr := ptr + 1;
16.            DFS_d(v);
17.        }
18.    }
19. }
```



```
1. // depth first search starting from vertex v of graph G
2. // input: starting node v
3. // output: f[|V|]: finish order
4. //          SCCs: resulting DFS forest
5. Algorithm DFS_d(v)
6. {
7.     visited[v] := 1;
8.     SCCs[ptr] := v; // store v in the same tree
9.     ptr := ptr + 1;
10.    for each vertex w adjacent to v do {
```

```

11.         if (visited[w] = 0) then DFS_d(w);
12.     }
13.     time := time + 1;
14.     f[v] := time; // record finish order
15. }

```

In depth first search, we traverse a graph in a manner that the if a vertex is unvisited, then we move to that vertex immediately, which makes the spanning tree deep. If we use adjacent matrix to store a graph, then the time complexity will be $\mathcal{O}(n^2)$ because in DFS_Call the loop run over 1 to n , and in DFS_d the loop must check n vertices to find all of its adjacent vertices. On the otherhand, if we use the adjacent list to store the graph, in DFS_Call we still need to run over 1 to n , but in DFS_d we already store the adjacent vertices together, therefore we don't need to run over all 1 to n to find its adjacent vertices. To be precise, we will spend a total of e times traversing the edges. Therefore, the time complexity is $\mathcal{O}(n + e)$

Besides storing the graph, both data structures need to store array f and $visited$, which takes $\mathcal{O}(n)$ additional space.

Time complexity: $\mathcal{O}(n + e)$ for adjacency list, $\mathcal{O}(n^2)$ for adjacency matrix

Overall space complexity: $\mathcal{O}(n) + \mathcal{O}(e)$ for adjacency list, $\mathcal{O}(n) + \mathcal{O}(n^2)$ for adjacency matrix

2.4 Sort Vertices

Because we are sorting elements based on their finish order, which are all integers, I think it would be a good opportunity to implement the counting sort.

Here is the generic counting sort:

```
1. // sort A[1 : n] in nondecreasing order
2. // and puts results into B[1 : n].
3. // assume  $1 \leq A[i] \leq k$ , for all  $i$ 
4. // input: A, int n, k
5. // output: B contains sorted results
6. Algorithm CountingSort_geric(A, B, n, k)
7. {
8.     // initialize C to all 0
9.     for  $i := 1$  to  $k$  do  $C[i] := 0$ ;
10.    // count # elements in C[A[i]]
11.    for  $i := 1$  to  $n$  do {
12.         $C[A[i]] := C[A[i]] + 1$ ;
13.    }
14.    // C[i] is the accumulate # of elements
15.    for  $i := 1$  to  $k$  do {
16.         $C[i] := C[i] + C[i - 1]$ ;
17.    }
18.    // store sorted order in array B
19.    for  $i := n$  to  $1$  step  $-1$  do {
20.         $B[C[A[i]]] := A[i]$ ;
21.         $C[A[i]] := C[A[i]] - 1$ ;
22.    }
23. }
```

Counting sort uses the value to be sorted, elements of array A, as index to store the position information in array C. When we need to store it to array B, we only need to look up for $A[i]$ in array C.

However, by just implement the above pseudo code will not get the job done.

In this case, we are sorting a list of indices with their finish time, rather than sorting the finish time itself. Therefore, we need some modification.

```
1. // sort idx[1 : n] in nonincreasing order using key
2. // assume  $1 \leq A[i] \leq k$ , for all  $i$ 
3. // input: idx, int n, k
4. // output: sorted idx
5. Algorithm CountingSort(idx, key, n, k)
6. {
7.     // initialize C to all 0
8.     for  $i := 1$  to  $k$  do  $C[i] := 0$ ;
9.     // count # elements in  $C[key[idx[i]]]$ 
10.    for  $i := 1$  to  $n$  do {
11.         $C[key[idx[i]]] := C[key[idx[i]]] + 1$ ;
12.    }
13.    //  $C[i]$  is the accumulate # of elements
14.    for  $i := 1$  to  $k$  do {
15.         $C[i] := C[i] + C[i - 1]$ ;
16.    }
17.    for  $i := n$  to  $1$  step  $-1$  do {
18.        // store sorted order back to array idx
19.         $\text{swap}(\text{idx}[n - C[key[idx[i]]]], \text{idx}[i])$ ;
20.         $C[key[idx[i]]] := C[key[idx[i]]] - 1$ ;
21.    }
22. }
```

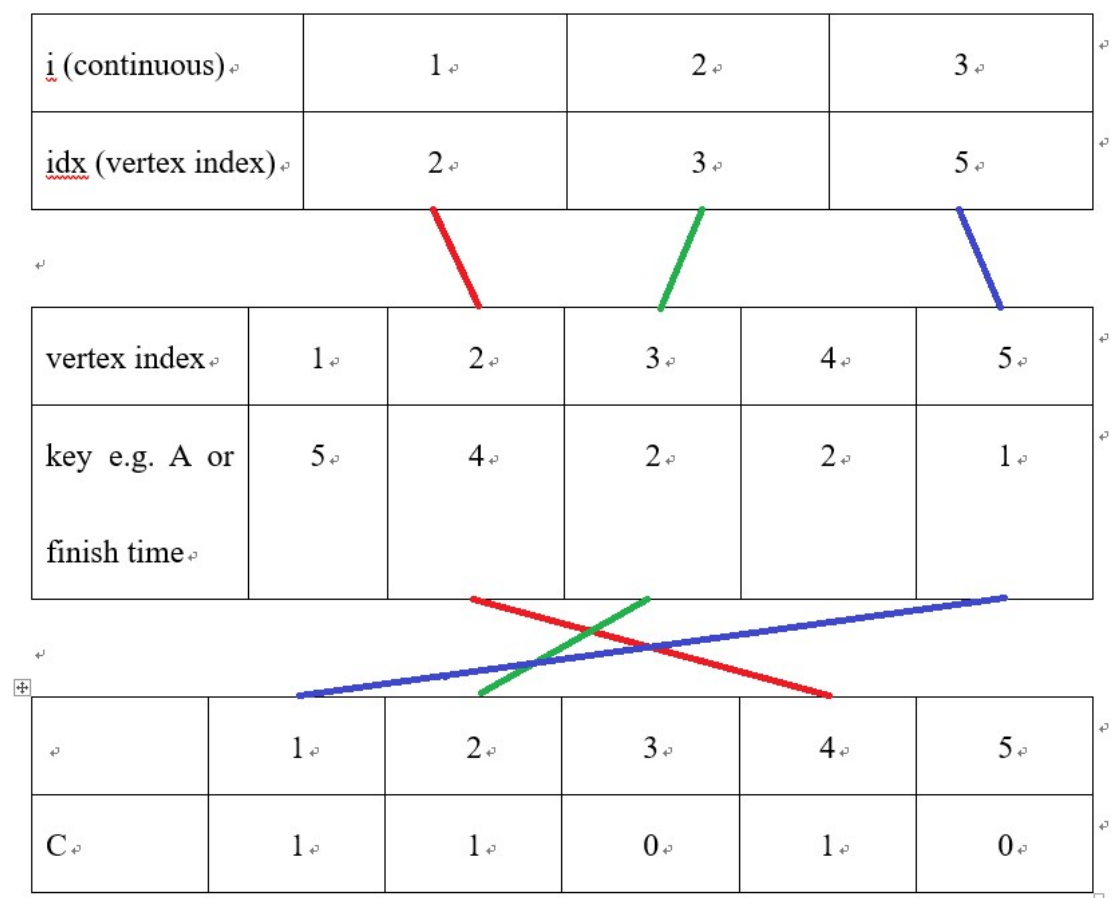
In `CountingSort_geric()`, we put the key $A[i]$ in array C using $A[i]$ as index, but in the modified version, we put $key[idx[i]]$ because key right now has the same roll as A and we no longer assume the index to be 1 to n continuously. idx is the mapping from 1 to n to the real index of the vertices. In the last loop, we do not need extra array B like we do in `CountingSort_geric()`, instead we

permute the elements of array `idx` to the correct place. Also, in line 19, we store the element backward to achieve the nonincreasing property.

For the time complexity, we are looping over 1 to n and k , therefore it is $\mathcal{O}(n+k)$. For the space complexity, besides the array to be sorted, we need additional array C , which takes up $\mathcal{O}(k)$ additional space.

Time complexity: $\mathcal{O}(n+k)$

Overall space complexity: $\mathcal{O}(n) + \mathcal{O}(k)$



Schematic diagram of CountingSort up to line 11

2.5 Summary of Implementation

As we can see, the analysis above depends heavily on how we store the graph. Because I want to strike a balance between space and time complexity.

I use dynamically allocated array to store the adjacent vertices. In the start, I initialize each vertex's the space to be e/n , and double the size if we run out of space. This way, the space complexity is better than adjacency matrix and a little worse than using linked list. We here assume each vertex has exactly e/n adjacent vertices, then the graph takes up $\mathcal{O}(e)$ space. If we analyze the total complexity:

1. construct the transposed graph G^T How?

Time complexity: $\mathcal{O}(e)$

Additional space complexity: $\mathcal{O}(e)$

2. DFS on G :

Time complexity: $\mathcal{O}(n + e)$

Additional space complexity: $\mathcal{O}(n)$

3. Sort vertices:

Calling CountingSort(e/n , n) n times

Time complexity: $n * \mathcal{O}(e/n + n) = \mathcal{O}(e + n^2) \approx \mathcal{O}(e)$

Space complexity: $\mathcal{O}(n)$

4. DFS on G^T

Time complexity: $\mathcal{O}(n + e)$

Additional space complexity: $\mathcal{O}(n)$

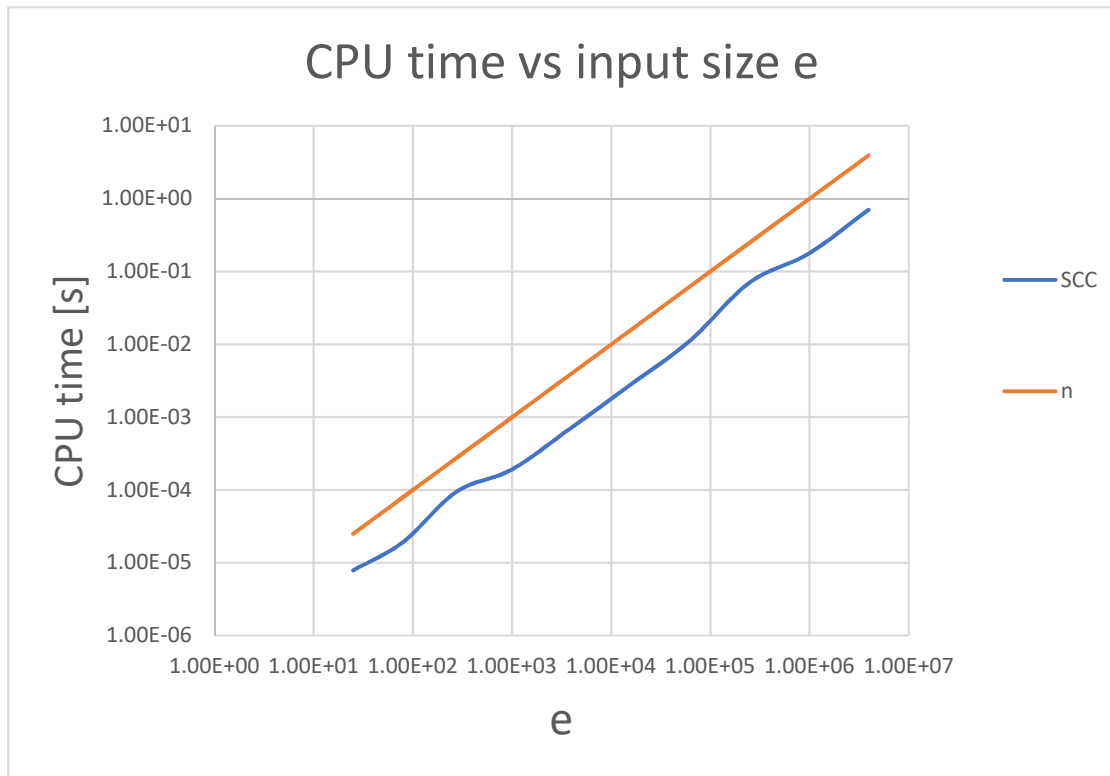
Therefore, the time total complexity is around $\mathcal{O}(e)$ and so does the space complexity.

3. Result and Observation

To measure the performance, we execute the program to solve the friend grouping problem with a different people and communication records.

e	CPU time [s]
25	7.87E-06
79	1.91E-05
274	9.39E-05
996	1.91E-04
3926	7.09E-04
15547	2.76E-03
61786	1.11E-02
246515	6.99E-02

984077	1.76E-01
3934372	7.02E-01



As we can see, the trend indeed follows the linear trend with input size e , thus validate our analysis.

hw07.c

```
1 // EE3980 HW07 Grouping Friends
2 // 105061110, 周柏宇
3 // 2020/04/24
4
5 #include <stdio.h>
6 #include <stdlib.h>
7 #include <string.h>
8 #include <sys/time.h>
9 #define N 100 // number of class in the hash table
10
11 // structure to store the index whose name is in the same class
12 typedef struct node {
13     int idx;
14     struct node *next;
15 } Node;
16
17 int n_names; // number of people
18 int n_links; // number of communication records
19 int n_SCCs; // number of friend groups
20 char **names; // array of names
21 int **adj_l; // adjacency list
22 int *adj_l_size; // size of list in adj_l
23 int *adj_l_ptr; // current position of list in adj_l
24 int **adj_lT; // adjacency list (transposed)
25 int *adj_lT_size; // size of list in adj_lT
26 int *adj_lT_ptr; // current position of list in adj_lT
27 int *f; // node's finish order
28 int *idx; // index mapping
29 int *visited; // record visited nodes
30 int time_DFS; // DFS clock
31 int *SCCs; // record nodes traversed by DFS
32 int SCCs_ptr; // current position of SCCs
33 int *C; // temporary array for CountingSort()
34 Node **bucket; // hash table
35
36 void readData(); // read data, construct the graph and its transposed version
37 void readData(void); // read data, construct the graph and its transposed versio
38 n
39 int nameToIndex(char *name); // convert name to array index
40 unsigned long hash(char *str); // hash function
41 double GetTime(void); // get local time in seconds
42 void SCC(); // find strongly connected components of a graph
43 void SCC(void); // find strongly connected components of a graph
44 void DFS_Call(int **G, int *len, int *idx);
45 // initialization and recursive DFS function call
46 void DFS_d(int **G, int *len, int v); // DFS from vertex v of the graph G
47 void CountingSort(int *idx, int *key, int n, int k);
48 // sort the idx by its key using counting sort
```

```

46 void freeAll(); // free all allocated memory
   void freeAll(void); // free all allocated memory
47
48 int main(void)
49 {
50     int i; // loop index
51     int component_idx = 0; // index of strongly connected component
52     double start, end; // timestamp
53
54     readData(); // read data, construct the graph and its transposed version
55     start = GetTime(); // start time
56     SCC(); // find strongly connected components of a graph
57     end = GetTime(); // end time
58     // print out input information, number of subgroups and execution time
59     printf("N = %d M = %d Subgroup = %d ", n_names, n_links, n_SCCs);
60     printf("CPU time = %.5e\n", end - start);
61     printf("Number of subgroups: %d\n", n_SCCs); // print out number of groups
62     for (i = 0; SCCs[i] != -2; i++) { // print out member of groups
63         if (i == 0) {
64             printf(" Subgroup %d:", ++component_idx);
65         }
66         else if (SCCs[i] == -1) {
67             printf("\n Subgroup %d:", ++component_idx);
68         }
69         else {
70             printf(" %s", names[SCCs[i]]);
71         }
72     }
73     printf("\n");
74     freeAll(); // free all allocated memory
75
76     return 0;
77 }
78
79 void readData()
   void readData(void)
80 {
81     int i, j, k, l; // indices
82     int class; // class of a name determined by the hash function
83     char tmp[20], tmp2[20], tmp3[20]; // temporary variable for input
84     int *newptr; // temporary pointer
85     Node *newNode; // temporary Node
86
87     // input number of people and communication records
88     scanf("%d %d", &n_names, &n_links);
89     // allocate memory
90     names = (char **)malloc(sizeof(char *) * n_names);
91     adj_l = (int **)malloc(sizeof(int *) * n_names);
92     adj_l_size = (int *)malloc(sizeof(int) * n_names);
93     adj_l_ptr = (int *)calloc(n_names, sizeof(int));

```

```

94 adj_lT = (int **)malloc(sizeof(int *) * n_names);
95 adj_lT_size = (int *)malloc(sizeof(int) * n_names);
96 adj_lT_ptr = (int *)calloc(n_names, sizeof(int));
97 bucket = (Node **)calloc(N, sizeof(Node *));
98 for (i = 0; i < n_names; i++) {
99     scanf("%s", tmp); // input names
100     names[i] = (char *)malloc(sizeof(char) * (3 * strlen(tmp) + 1));
101     strcpy(names[i], tmp);
102     class = hash(tmp) % N; // decide the name's class
103     if (!bucket[class]) { // start of the linked list
104         bucket[class] = (Node *)malloc(sizeof(Node));
105         bucket[class]->idx = i; // store the index at this class
106         bucket[class]->next = NULL;
107     }
108     else { // this class already has data
109         newNode = (Node *)malloc(sizeof(Node));
110         newNode->idx = i; // store the index at this class
111         newNode->next = bucket[class];
112         bucket[class] = newNode;
113     }
114 }
115 for (i = 0; i < n_names; i++) {
116     // initialize the adjacency list for the graph
117     adj_l_size[i] = n_links / n_names;
118     adj_l[i] = (int *)malloc(sizeof(int) * adj_l_size[i]);
119     // initialize the adjacency list for the transposed graph
120     adj_lT_size[i] = n_links / n_names;
121     adj_lT[i] = (int *)malloc(sizeof(int) * adj_lT_size[i]);
122 }
123 for (i = 0; i < n_links; i++) {
124     // input communication records
125     scanf("%s %s %s", tmp, tmp2, tmp3);
126     // convert names to array indices
127     j = nameToIndex(tmp);
128     k = nameToIndex(tmp3);
129
130     // dynamically allocate memory for the adjacency list
131     if (adj_l_ptr[j] >= adj_l_size[j]) {
132         adj_l_size[j] *= 2; // double the size
133         newptr = (int *)malloc(sizeof(int) * adj_l_size[j]);
134         // allocate memory with new size
135         for (l = 0; l < adj_l_ptr[j]; l++) newptr[l] = adj_l[j][l];
136         // copy the data to new memory
137         free(adj_l[j]); // free old memory block
138         adj_l[j] = newptr; // assign the new pointer
139     }
140     adj_l[j][adj_l_ptr[j]] = k; // store the edge <j, k>
141     adj_l_ptr[j]++; // update current position in the list
142
143     // dynamically allocate memory for the adjacency list

```

```

144     if (adj_lT_ptr[k] >= adj_lT_size[k]) {
145         adj_lT_size[k] *= 2; // double the size
146         newptr = (int *)malloc(sizeof(int) * adj_lT_size[k]);
147                                     // allocate memory with new size
148         for (l = 0; l < adj_lT_ptr[k]; l++) newptr[l] = adj_lT[k][l];
149                                     // copy the data to new memory
150         free(adj_lT[k]); // free old memory block
151         adj_lT[k] = newptr; // assign the new pointer
152     }
153     adj_lT[k][adj_lT_ptr[k]] = j; // store the edge <k, j>
154     adj_lT_ptr[k]++; // update current position in the list
155 }
156 }
157
158 int nameToIndex(char *name) // convert name to array index
159 {
160     Node *tmp; // temporary Node
161
162     tmp = bucket[hash(name) % N];
163     while (tmp) { // return the index if the names are matched
164         if (!strcmp(name, names[tmp->idx])) return tmp->idx;
165         tmp = tmp->next;
166     }
167
168     return -1;
169 }
170
171 unsigned long hash(char *str) // hash function
172 {
173     unsigned long hash = 5381;
174     int c;
175
176     while ((c = *str++)) {
177         hash = ((hash << 5) + hash) + c;
178     }
179
180     return hash;
181 }
182
183 double GetTime(void) // get local time in seconds
184 {
185     struct timeval tv; // variable to store time
186
187     gettimeofday(&tv, NULL); // get local time
188
189     return tv.tv_sec + 1e-6 * tv.tv_usec; // return local time in seconds
190 }
191
192 void SCC() // find strongly connected components of a graph
193     void SCC(void) // find strongly connected components of a graph

```

```

193 {
194     int i; // loop index
195
196     // allocate memory
197     f = (int *)malloc(sizeof(int) * n_names);
198     idx = (int *)malloc(sizeof(int) * n_names);
199     SCCs = (int *)calloc(n_names * 2 + 1, sizeof(int));
200     visited = (int *)malloc(sizeof(int) * n_names);
201     C = (int *)malloc(sizeof(int) * n_names);
202
203     // initialize the index
204     for (i = 0; i < n_names; i++) idx[i] = i;
205
206     // traverse the graph in the given order
207     DFS_Call(adj_l, adj_l_ptr, idx);
208
209     // sort the index using the array f as key (decreasing order)
210     CountingSort(idx, f, n_names, n_names);
211
212     // sort the adjacency list of the transposed graph in-place
213     // using the array f as key (decreasing order)
214     for (i = 0; i < n_names; i++) {
215         CountingSort(adj_lT[i], f, adj_lT_ptr[i] - 1, n_names);
216     }
217
218     // traverse the transposed graph in the given order
219     DFS_Call(adj_lT, adj_lT_ptr, idx);
220 }
221
222 // initialization and recursive DFS function call
223 void DFS_Call(int **G, int *len, int *idx)
224 {
225     int i, j; // loop index
226
227     // initialize
228     for (i = 0; i < n_names; i++) visited[i] = 0;
229     time_DFS = 0;
230     SCCs_ptr = 0;
231     n_SCCs = 0;
232     for (i = 0; i < n_names; i++) {
233         j = idx[i]; // decide the vertex to travel according idx
234         if (visited[j] == 0) {
235             SCCs[SCCs_ptr++] = -1; // separate different subgroups with -1
236             n_SCCs++; // number of subgroups increase by one
237             DFS_d(G, len, j); // start DFS from vertex j
238         }
239     }
240     SCCs[SCCs_ptr] = -2; // end of array
241 }
242

```

```

243
244 void DFS_d(int **G, int *len, int v) // DFS from vertex v of the graph G
245 {
246     int i, j; // loop index
247
248     visited[v] = 1; // vertex v has been visited
249     SCCs[SCCs_ptr++] = v; // add vertex v to the current subgroup
250     for (i = 0; i < len[v]; i++) {
251         j = G[v][i]; // decide next vertex to travel
252         if (visited[j] == 0) {
253             DFS_d(G, len, j); // DFS from vertex j
254         }
255     }
256     f[v] = time_DFS++; // record the finishing order
257 }
258
259 // sort the idx by its key using counting sort
260 void CountingSort(int *idx, int *key, int n, int k)
261 {
262     int i; // loop index
263     int tmp, tmp_idx; // temporary variable
264
265     for (i = 0; i < k; i++) C[i] = 0; // initialize C to all 0
266     for (i = 0; i < n; i++) {
267         C[key[idx[i]]]++; // count # elements in C[key[idx[i]]]
268     }
269     for (i = 1; i < k; i++) {
270         C[i] += C[i - 1]; // C[i] is the accumulate # of elements
271     }
272     for (i = n - 1; i >= 0; i--) {
273         // store sorted order back to array idx
274         tmp_idx = n - (C[key[idx[i]]] - 1) - 1; // decreasing order
275         tmp = idx[tmp_idx];
276         idx[tmp_idx] = idx[i];
277         idx[i] = tmp;
278         C[key[idx[i]]]--;
279     }
280 }
281
282 void freeAll() // free all allocated memory
283 void freeAll(void) // free all allocated memory
284 {
285     int i; // loop index
286     Node *tmp, *next;
287
288     for (i = 0; i < n_names; i++) {
289         free(names[i]);
290         free(adj_l[i]);
291         free(adj_lT[i]);
292     }

```



```

292     for (i = 0; i < N; i++) {
293         tmp = bucket[i];
294         while (tmp) {
295             next = tmp->next;
296             free(tmp);
297             tmp = next;
298         }
299     }
300     free(names);
301     free(adj_l);
302     free(adj_l_size);
303     free(adj_l_ptr);
304     free(adj_lT);
305     free(adj_lT_size);
306     free(adj_lT_ptr);
307     free(visited);
308     free(f);
309     free(idx);
310     free(SCCs);
311     free(C);
312     free(bucket);
313 }

```

[Program Format] can be improved.

[Writing] hw07a.pdf spelling errors: vertexs(1)

[CPU] time for c10.dat: 1.47000e-01 sec

[Approach] counting sort is not needed, DFS can be used to get sorted results.

[Approach] what is the transposed graph and how did you create one?

[Report] writing can be more clear.

Score: 81