

**Problem 1: Graph cohesiveness (KT 7.46)**

In sociology, one often studies a graph  $G$  in which nodes represent people and edges represent those who are friends with each other. Let's assume for purposes of this question that friendship is symmetric, so we can consider an undirected graph.

Now suppose we want to study this graph  $G$ , looking for a "close-knit" group of people. One way to formalize this notion would be as follows. For a non-empty subset  $S$  of nodes, let  $e(S)$  denote the number of edges in  $S$ -that is, the number of edges that have both ends in  $S$ . We define the *cohesiveness* of  $S$  as  $e(S)/|S|$ . A natural thing to search for would be a set  $S$  of people achieving the maximum cohesiveness.

1. Give a polynomial-time algorithm that takes a rational number  $\alpha$  and determines whether there exists a set  $S$  with cohesiveness greater than  $\alpha$ .
2. Give a polynomial-time algorithm to find a set  $S$  of nodes with maximum cohesiveness.

Problem 1.1Goal:

Determine whether  $\exists$  a "close-knit" group of people in  $G$

s.t.

$$\frac{\text{the # of edges w/ both endpoints in a subset of nodes } (S) \text{ in } G}{\text{The # of nodes in } S} > \alpha$$

Definitions:

$S$  = a non-empty subset of nodes in  $G$

$|S|$  = The number of nodes in  $S$

$e(S)$  = The number of edges in  $S$  (that have both their endpoints in  $S$ )

$$\text{Cohesiveness} = \frac{e(S)}{|S|}$$

$\alpha$  = rational number for which want to check if  $\exists$  a set  $S$  in  $G$  with cohesiveness greater than this

- As  $\alpha$  is a rational number we can rewrite it as an integer fraction  $\rightarrow \alpha = \frac{a}{b}$

$a$  = Numerator for the reformulated  $\alpha$  parameter

$b$  = Denominator for the reformulated  $\alpha$  parameter.

### Key Idea:

Given an undirected friend graph  $G$ , we can create an equivalent bipartite flow graph  $G'$  with edges  $E'$  and vertices  $V$  from  $G$  as vertices in  $G'$ . Moreover, this can be reformulated as a flow graph by adding a source node  $s$  and a sink node  $t$ .

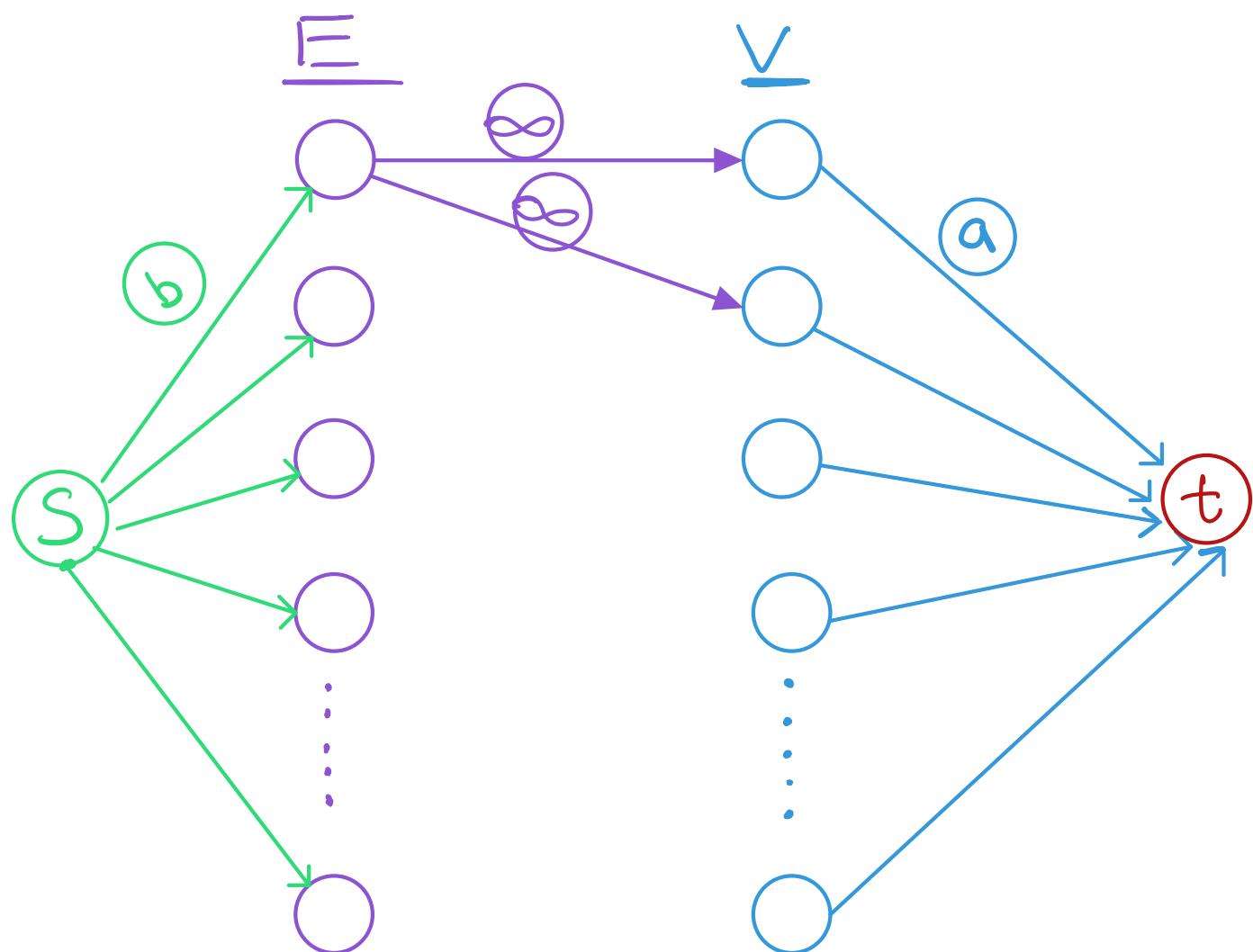
To create a problem equivalency (i.e., valid reduction) we can use the  $\alpha$  formulation defined above to check if cohesiveness. The existence of  $\alpha$  cohesiveness will relate to the min cut/max flow, and the exact subset  $S$  yielding  $>\alpha$  cohesiveness (if one exists) can be extracted from the residual graph  $R_G$  after the completion of the max flow algorithm.

Specifically the quantity that is being evaluated / maximized in the context of this problem is:

$$\frac{e(s)}{\|S\|} > \alpha = \frac{a}{b} \Rightarrow b \cdot e(s) - a \cdot \|S\| > 0$$

Thus if there exist a  $b \cdot e(s) - a \cdot \|S\| > 0$ , then there exists a set  $S$  in  $G$  with a cohesiveness  $> \alpha$ .

### Graph Construction:



Define the above flow graph  $G'$  with the following:

$V \rightarrow$  the set of vertices in  $G'$

$E \rightarrow$  the set of edges in  $G'$

$V = \{s, t\} \cup \{V_i \mid \text{vertices } V_i \text{ in } G\}$   
 $\cup \{V \mid \text{undirected edges } e_{ij} \text{ in } G\}$

$E = \{(s, e_{ij}) \mid \begin{array}{l} \text{edges between the source and edge vertices} \\ \text{in } G' \end{array}\}$

with a capacity of  $b$ .

$\cup \{(e_{ij}, V_i), (e_{ij}, V_j) \mid i \neq j\}$  with a conceptual  
capacity of  $\infty$ .

$\cup \{(V_i, t) \mid \begin{array}{l} \text{edges between all } G \text{ vertices and the sink in} \\ G' \end{array}\}$  with a capacity of  $a$ .

## Algorithm:

Create a bipartite flow graph  $G'$  from  $G$  as described above in Graph Construction. Run one of the four max flow algorithms (for selection see Time Complexity and Max Flow Algorithm Selection below). If the max flow/mincut is  $< b \cdot \text{the\#}$  of edges from  $S$  to  $E$ , conclude that there exists a set  $S$  in  $G$  with a cohesiveness  $> \alpha$ . The valid set  $S$  with cohesiveness  $> \alpha$  can be extracted by finding all nodes reachable from  $S$  in  $R_{G'}$  at the termination of the max flow algorithm (i.e., there are no augmenting paths in  $R_{G'}$  from  $S$  to  $t$ ). Specifically this can be extracted by running breadth-first search (BFS) on  $R_{G'}$  from  $S$  and returning all reachable vertices in  $V$  and edges in  $E$ . In the event the max flow/mincut is  $\geq b \cdot \text{the\#}$  of edges from  $S$  to  $E$ , return that there does not exist a set  $S$  in  $G$  with a cohesiveness  $> \alpha$ .

Implementation Intricacy: The maximum flow can be checked by summing the flow out of the source node  $S$ , along its outgoing edges. Also note the # of edges from  $S$  to  $E$  are static in  $G'$  & can be computed easily as this is the bonding capacity out of  $S$ . Thus  $b \cdot$  bonding capacity of  $S$  can be computed before running the maxflow algorithm & can be stored to compare against the calculated max flow.

## Proof of Correctness:

### CLAIM 1:

If the minimum cut (/max flow) is  $\sum_{i \in S} b_i$  if from  $S$  to  $E$   
then  $\exists$  a set  $S$  in  $G$  with cohesiveness  $> \alpha$ .

### CLAIM 2:

The minimum cut must necessarily omit conceptual  $\infty$  capacity edges, and thus the cut capacity will be determined by edges from  $S$  to  $E$  and  $V$  to  $t$ .

### CLAIM 3:

If there exists a set  $S$  in  $G$  with cohesiveness  $> \alpha$ , then the set of reachable edges in  $R_{G'}$  from  $S$  corresponds to a valid  $S$  satisfying the problem constraints.

- The validity of CLAIM 2 can be easily shown. Specifically there always exists a finite capacity cut that omits  $\infty$  capacity edges (such as the cut of  $S$  from the rest of  $G'$  or  $t$  from the rest of  $G'$ ). Thus no infinite edges can be in a min cut.
- CLAIM 1 and CLAIM 3 are related and relate to the problem reduction to a flow graph, and the extraction of the problem answer from running the max flow/min cut algorithm corresponds to the problem answer.

- Let  $A$  be the set of vertices in  $RG'$  s.t. they can be reached from the source node  $S$ .
- Let  $B$  be the complement of  $A$ , that is  $B = V - A$
- Know that at the completion of max flow algorithm  $CAP(A, B) = \text{maxflow}$

$$CAP(A, B) = \sum_{i \notin A} b + \sum_{j \in A} a$$

Define  $K$  to be  $\sum_{i \notin A} b$  if from  $S$  to  $E$  → This is the cut capacity of  $A = S$  and thus is a constant value

Observe that we can manipulate  $CAP(A, B)$  by adding &

subtracting  $\sum_{i \in A} b$ :

$$CAP(A, B) = \underbrace{\sum_{i \notin A} b + \sum_{i \in A} b}_{=K} - \sum_{i \in A} b + \sum_{j \in A} a$$

$$CAP(A, B) = K - \sum_{i \in A} b + \sum_{j \in A} a$$

$$= K - \left[ \sum_{i \in A} b - \sum_{j \in A} a \right]$$

Note that by the construction of  $G'$  that

$$\sum_{i \in A} b - \sum_{j \in A} a = b \cdot e(s) - a \cdot \|s\|$$

as all edges from  $E$  to  $V$  in  $G$  were maintained and  
Capacity unbounded (i.e., conceptual capacities of  $\infty$ )

Thus here we have:

$$CAP(A, B) = b - [b \cdot e(s) - a \cdot \|s\|]$$

- The minimum cut then occurs when the problem objective is maximized!
- Thus if the min cut  $< b$ , then there exists a set  $S$  in  $G$  because necessarily  $b \cdot e(s) - a \cdot \|s\| > 0$ !
- Finally by definition of  $R_{G'}$  the set of all reachable vertices from  $s$  in this case corresponds to the edges and vertices in  $G$  that exist in  $S$  (excluding the source node  $s$ ).

## Time Complexity and Max Flow Algorithm Selection:

Creation of  $G'$  and the individual steps of the max flow algorithm are linear in the size of  $G'$ . Running BFS on  $R_{G'}$  is also linear in the size of  $G'$ . Thus the time complexity bottleneck is the max flow algorithm selected.

Given  $G'$  as defined in Graph Construction observe the following size of  $G'$ :

$$\|V\| = \|E\| + \|V\| + 2$$

Define  $\|V\| = n$   
for clarity

$$\|E\| = \underbrace{\|E\|}_{\text{layer 1 edges}} + \underbrace{2 \cdot \|E\|}_{\text{layer 2 edges}} + \underbrace{\|V\|}_{\substack{\text{layer 3 } (v \rightarrow t) \\ \text{edges}}}$$

Note in the worst case  $G$  is fully connected and thus there are  $\binom{n}{2}$  undirected edges in  $G$ . This then gives

$$\|V\| = \frac{n(n-1)}{2} + n + 2 \approx n^2$$

$$\|E\| = \frac{n(n-1)}{2} + n \cdot (n-1) + n \approx n^2$$

Using the highest order terms for both  $\|V\|$  and  $\|E\|$  set the following time complexities for each of the max flow algorithms.

Ford-Fulkerson:  $(|V| + |E|) \cdot \text{Max flow}$  (or  $G_{\text{total}}$  if unknown)

- Here this equates to  $2n^2 \cdot \text{Max flow}$
- In the worst case we know max flow is bounded by  $K$  (defined above  $\rightarrow$  occurs when there does not exist an  $S$  in  $G$  with cohesiveness  $> 2$ ).

• Thus in the worst case this approach yields a runtime of  $2n^2 K$

Scaling:  $(|V| + |E|)^2 \cdot \log C_{\max}$

- Here this equates to  $4n^4 \cdot \log(\max(b, a))$

Barg-edmonds-dinitz:  $(|V| + |E|)^2 \cdot |V|$

- Here this equates to  $4n^6$

Preflow Push:  $(|V|)^3$

- Here this equates to  $n^6$

Given this we can discard Barg-edmonds-dinitz as it is always worse than Preflow Push by a constant factor. Thus the choice of algorithm is dependent on the connectivity of  $G$  and the given  $\alpha$  parameter.

Know there is an upper bound to  $\alpha$  for feasibility as

$$\frac{\binom{n}{2}}{n} = \frac{a}{b} \Rightarrow \frac{(n-1)}{2} \approx n$$

Thus we can bound scaling and Ford-Fulkerson in terms of  $n$

$$\text{Scaling} \rightarrow 4n^4 \cdot \log(\max(b, a)) \quad \leftarrow a \text{ is the max in worst case w/ } \binom{n}{2}$$

$$\begin{aligned} &\cdot \text{Goes to } 4n^4 \cdot \log n^2 \\ &= 8n^4 \log n \end{aligned}$$

$$\text{Ford-Fulkerson} \rightarrow 2n^2 K$$

$$K = \sum_{\text{if from } S \text{ to } E} b \quad \text{with } b \sim n$$

• Recall in the worst case there are  $\binom{n}{2}$  edges in  $G$

$$\cdot \text{Thus } K \approx n^3$$

$$\cdot \text{This gives: } 2n^2(n^3) = 2n^5$$

Given these time complexities, observe Scaling has the lowest order run time of the above 4 approaches. Thus select Scaling and conclude a time complexity of  $\mathcal{O}(8n^4 \log n) \rightarrow \mathcal{O}(n^4 \log n)$ .

## Problem 1.2:

Key Idea:

We can use the approach developed in Problem 1.1 above to find the maximum cohesiveness in  $G$ . Specifically we can perform a binary search across  $\alpha$  values. There are between  $0$  and  $\binom{n}{2}$  possibilities for  $e(S)$  and  $n$  possibilities for  $S$  yielding a search space of  $n^3$ .

Algorithm:

Start from the  $\alpha$  midpoint (i.e.,  $\frac{\binom{n}{2}}{2} = \frac{n-1}{4}$ )

and run the algorithm described above in Problem 1.1. If there exists an  $S$  with that cohesiveness, store the current  $S$  and go "right" by changing  $\alpha$  to the midpoint between my current  $\alpha$  and the current max right (which is initialized at the upper bound  $\frac{n-1}{2}$ ). Set the min left here to be the previously searched  $\alpha$ . If there does not exist an  $S$  with cohesiveness  $\alpha$ , go "left" by changing  $\alpha$  to the midpoint between my current  $\alpha$  and the min left (initialized to 0). Set the max right to my previously searched  $\alpha$ .

Continue this process until  $\min \text{left} = \max \text{right}$  and return the stored  $S$  (or null in the case of a completely disconnected graph  $G$ ).

### Time Complexity:

Performing binary search across  $n^3$  values gives a time complexity of  $O(\log n^3) \rightarrow O(3 \log n) \rightarrow O(\log n)$ . Combining this with the Scaling Selection algorithm selected for the algorithm in Problem 1.1 yields the following time complexity:

$$O((8n^4 \cdot \log n) \cdot 3 \log n) \rightarrow O(24n^4 \log^2 n)$$

Removing constant factors yields the final time complexity of  $O(n^4 \log^2 n)$ .

### Proof of Correctness:

Known by the proof outlined for Problem 1.1 above that for a particular of the approach is valid. Here observe that binary search allows us to narrow the search space. Specifically the key point is that if there exists no set  $S$  with a given  $d$ , then there can exist no set  $S$  with an even greater  $d$ . Thus we can validly narrow the search space by a factor of 2 at each iteration.

### Problem 2: Remote Sensors

Devise as efficient as possible algorithm for the following problem. You have  $n$  remote sensors  $s_i$  and  $m < n$  base stations  $B_j$ . For  $1 \leq j \leq m$ , base station  $B_j$  is located at  $(x_j, y_j)$  in the two-dimensional plane. You are given that no two base-stations are less than 1 km apart (in standard Euclidean distance,  $\sqrt{((x_j - x_k)^2 + (y_j - y_k)^2)}$ ). All base stations have the same integer bandwidth capacity  $C$ .

For  $1 \leq i \leq n$ , sensor  $s_i$  is located at  $(x_i, y_i)$  in the two-dimensional plane and has an integer bandwidth requirement of  $r_i$ , which can be met by assigning bandwidth on multiple base stations. Let  $b_{i,j}$  be the amount of bandwidth assigned to sensor  $s_i$  on base station  $B_j$ . The assignment must meet the following constraints:

- No sensor may be assigned any bandwidth on a base station more than 2 km distance from it, i.e., if the distance from  $s_i$  to  $B_j$  is greater than 2,  $b_{i,j} = 0$ .
- The sum of all the bandwidth assigned to any remote sensor  $s_i$  must be at least  $r_i$ : for each  $1 \leq i \leq n$ ,  $\sum_j b_{i,j} \geq r_i$ .
- The sum of all bandwidth assigned on base station  $B_j$  must be at most  $C$ : for each  $1 \leq j \leq m$ ,  $\sum_i b_{i,j} \leq C$ .

Your algorithm should find a solution meeting the above constraints if possible, and otherwise output a message saying "No solution exists". Prove the correctness of your algorithm and discuss its time complexity.

### Definitions:

Let  $m$  = the number of base stations

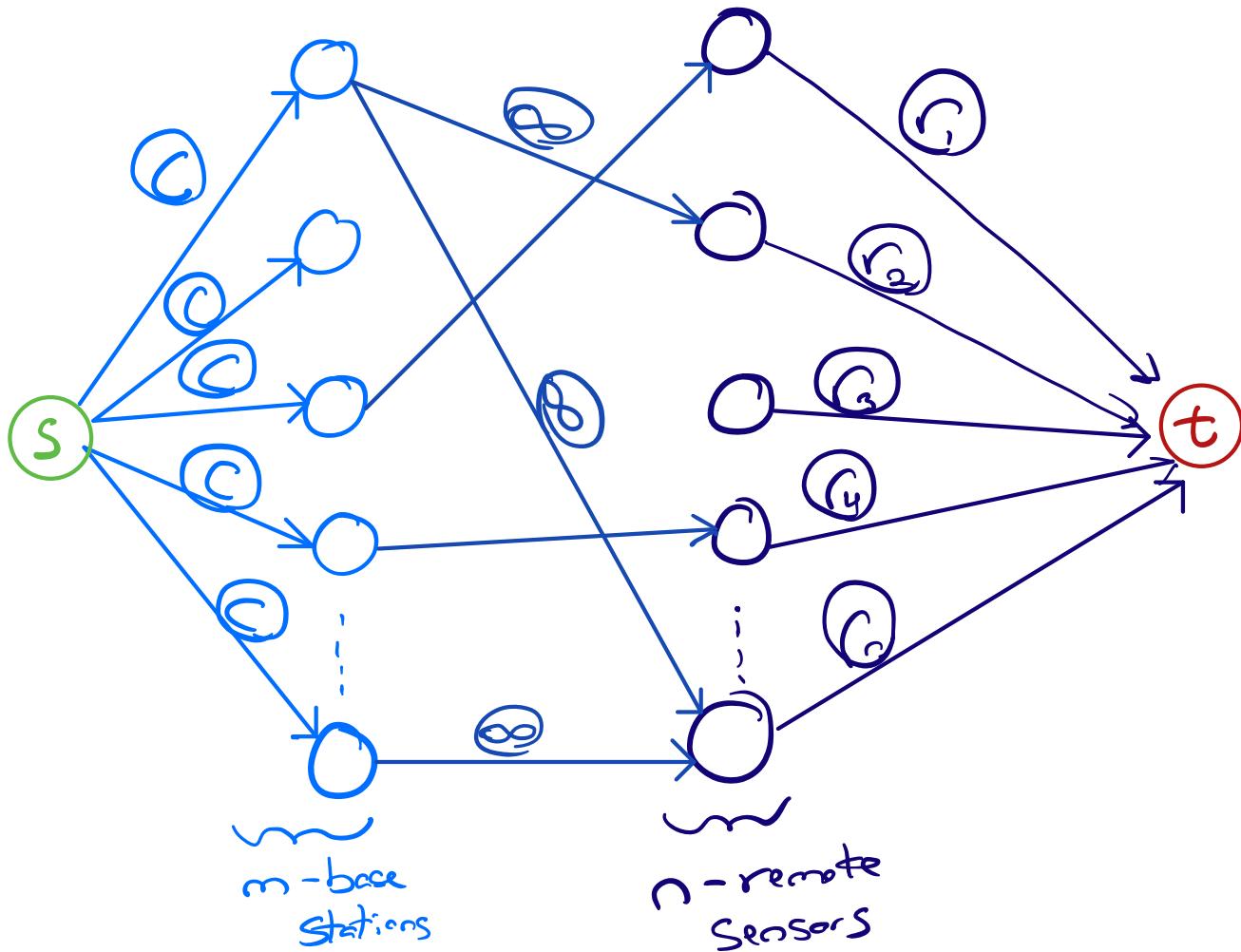
Let  $n$  = the number of remote sensors

Let  $C$  = the integer bandwidth capacity of a base station

### Key Idea:

Each sensor has an integer bandwidth requirement  $r_i$ , which can be satisfied by assigning bandwidth ( $< C$ ) from all sensors within a valid radius (i.e., 2 Km). The radius of base stations, coupled with the assignment nature of the problem, lends itself to the formulation of a bipartite graph, and maximizing flow from base stations to sensors.

## Graph Construction:



Define the above graph  $G$  with the following:

The set of vertices  $V$ :

$$V = \{s\} \cup \{B_j \mid 1 \leq j \leq m\} \cup \{S_i \mid 1 \leq i \leq n\} \cup \{t\}$$

Edges E:

$E = \{(S, B_j) \mid 1 \leq j \leq m\}$  with capacity of C  
the bandwidth integer capacity of the base stations.

$\cup \{(B_i, S_i) \mid \text{if } B_i \text{ w/i } 2\text{km of a sensor } S_i\}$   
with a conceptual capacity of  $\infty$ .

$\cup \{(S_i, t) \mid 1 \leq i \leq n\}$  with a capacity of  $r_i$  the  
integer bandwidth requirement for sensor  $i$

### Algorithm:

Create a bipartite graph  $G$  as defined above in Graph Construction. Run the Ford-Fulkerson algorithm on  $G$  to obtain the maximum flow. If the maximum flow equals  $\sum_{i=1}^n r_i$  (i.e., the bandwidth requirements for each remote sensor  $S_i$ ) then the solution of all  $b_{ij}$  that satisfies the problem constraints is the flow along the  $(B_j, S_i)$  edges. If the max flow is  $< \sum_{i=1}^n r_i$  then return the message "No solution exists" as specified by the problem. Implementation Intrinsic: The maximum flow can be checked by summing the flow out of the source node  $S$ , along its outgoing edges.

## Proof of Correctness:

### CLAIM 1:

$\exists$  an assignment of bandwidth from stations to sensors satisfying the problem constraints if  $\text{max flow} = \sum_{i=1}^n r_i$ , otherwise no valid assignment exists.

- Note the capacities of the edges  $(S, B_j)$  and  $(S_i, t)$ . Satisfy the upper & lower bounds of the problem constraints. Specifically, a base station can push no more than  $r_i$  bandwidth and a sensor must receive at least  $r_i$  bandwidth from base stations w/i 2 Km from itself.
- These constraints are satisfied by the edge capacities in  $G$  + the  $(B_j, S_i)$  edges which are only constructed between base stations and a sensor w/i 2 Km of one another.
- Note that when  $\text{max flow} < \sum_{i=1}^n r_i$ , it means that the limitations occur either through bandwidth upper limits of a base station or the lack of an edge to push bandwidth to a sensor (i.e., a base station is "out of range").
- Also note that  $\text{max flow} = \sum_{i=1}^n r_i$  cannot occur erroneously. By setting the lower bound of a sensors bandwidth requirement as the edge capacity from  $(S_i, t)$  this ensures no

superfluous flow  $> r_i$  can occur for all  $i$  demands.  
 Thus only when the max flow =  $\sum_{i=1}^n r_i$  can an assignment of  
 bandwidth from base stations to sensors occur.

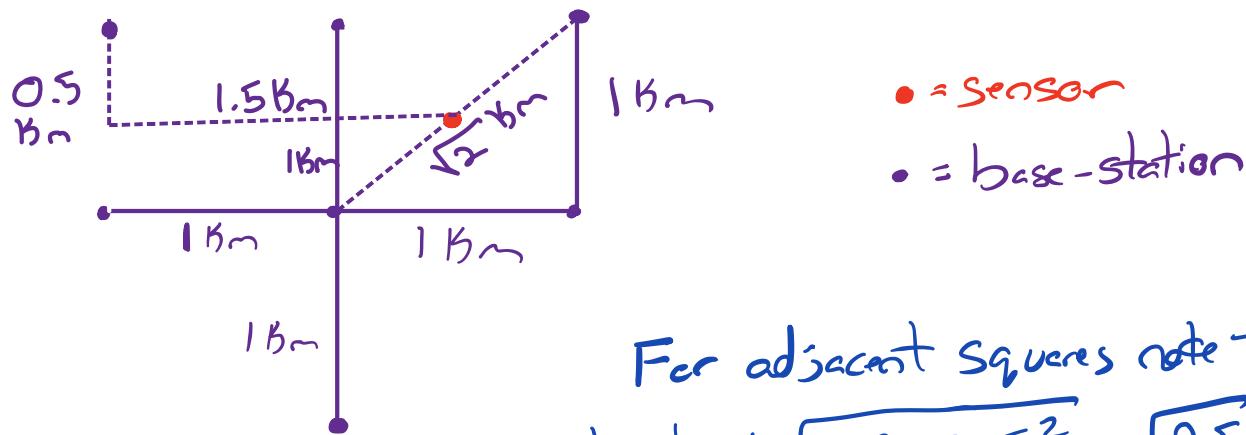
### CLAIM 2:

If the max flow =  $\sum_{i=1}^n r_i$ , then the flow along  $(B_j, S_i)$  edges  
 corresponds to the bandwidth assignment  $b_{ij}$  from base  
 stations to sensors.

- The flow bottleneck for a base station  $B_j$  is  $C$  & given the construction of  $G$ , all the sensors to which  $B_j$  is connected can have at most flow =  $C$ .
- All the sensors to which  $B_j$  is connected then can have a flow that is  $\leq \min(C, r_i)$  given the capacities of the edges.
- Edges  $(B_j, S_i)$  will only have flow if  $\exists$  an edge (i.e.,  $B_j$  is within  $2K_m$  of  $S_i$ ) and the Ford-Fulkerson algorithm assigns flow to an edge.
- Thus by the edge capacities, the construction of edges  $(B_j, S_i)$  and CLAIM 1, the flow along these edges corresponds to a valid assignment.

## Time Complexity and Max Flow Algorithm Selection:

Note that the problem states that "no two base-stations are less than 1 Km apart". This equates to saying that for any base station,  $\exists$  a radius of 1 w/i which no base-station can also exist. In the worst case scenario then we observe that base-stations could establish a lattice-like pattern, with a sensor placed equidistant between all sensors in a "square".



For adjacent squares note the distance equates to:  $\sqrt{1.5^2 + 0.5^2} = \sqrt{2.5} \text{ Km} < 2 \text{ Km}$

For subsequent corners, observe the distance is  $\frac{3\sqrt{2}}{2} > 2 \text{ Km}$

Thus in the worst case a sensor can have an edge between itself and 12 base-stations.

12 occurs from the 4 corners of my current square + 2 corners in the left, right, above, and below squares.

Given this we know that  $|V| = m+n+2$  and

$$|E| = m+n+12n = m+13n \text{ in the worst case.}$$

• As we know the max flow is bounded by  $\sum_{i=1}^n r_i$ , we observe

Ford-Fulkerson has a time complexity of:

$$(|V| + |E|) \cdot \text{Max Flow} = (2m + 14n + 2) \cdot \sum_{i=1}^n r_i$$

Preflow-Push, Scaling, and KED are all polynomial in  $m$  and  $n$  & thus select the Ford-Fulkerson algorithm for  $\sum_{i=1}^n r_i$

$$< \min((m+n+2)^3, (2m+14n+2) \cdot \min(\log C_{\max}, m+n+2)).$$

Thus conclude that the time complexity of this algorithm is  $\mathcal{O}((2m+14n+2) \cdot \sum_{i=1}^n r_i) \rightarrow \mathcal{O}((m+n) \sum_{i=1}^n r_i)$  by selection of Ford-Fulkerson.

Note:

Graph creation & individual steps of the algorithm are linear in the size of the graph, & thus conclude the flow algorithm above is the time complexity bottleneck.

### **Problem 3: Scheduling in a medical consulting firm (KT 7.19)**

You've periodically helped the medical consulting firm Doctors Without Weekends on various hospital scheduling issues, and they've just come to you with a new problem. For each of the next  $n$  days, the hospital has determined the number of doctors they want on hand; thus, on day  $i$ , they have a requirement that exactly  $p_i$  doctors be present.

There are  $k$  doctors, and each is asked to provide a list of days on which he or she is willing to work. Thus doctor  $j$  provides a set  $L_j$  of days on which he or she is willing to work.

The system produced by the consulting firm should take these lists and try to return to each doctor  $j$  a list  $L'_j$  with the following properties.

- (A)  $L'_j$  is a subset of  $L_j$ , so that doctor  $j$  only works on days he or she finds acceptable.
- (B) If we consider the whole set of lists  $L'_1, \dots, L'_k$ , it causes exactly  $p_i$  doctors to be present on day  $i$ , for  $i = 1, 2, \dots, n$ .

1. Describe a polynomial-time algorithm that implements this system. Specifically, give a polynomial-time algorithm that takes the numbers  $p_1, p_2, \dots, p_n$ , and the lists  $L_1, \dots, L_k$ , and does one of the following two things.

- Return lists  $L'_1, L'_2, \dots, L'_k$  satisfying properties (A) and (B); or
- Report (correctly) that there is no set of lists  $L'_1, L'_2, \dots, L'_k$  that satisfies both properties (A) and (B).

2. The hospital finds that the doctors tend to submit lists that are much too restrictive, and so it often happens that the system reports (correctly, but unfortunately) that no acceptable set of lists  $L'_1, L'_2, \dots, L'_k$  exists.

Thus the hospital relaxes the requirements as follows. They add a new parameter  $c > 0$ , and the system now should try to return to each doctor  $j$  a list  $L'_j$  with the following properties.

(A\*)  $L'_j$  contains at most  $c$  days that do not appear on the list  $L_j$ .

(B) (Same as before) If we consider the whole set of lists  $L'_1, \dots, L'_k$ , it causes exactly  $p_i$  doctors to be present on day  $i$ , for  $i = 1, 2, \dots, n$ .

Describe a polynomial-time algorithm that implements this revised system. It should take the numbers  $p_1, p_2, \dots, p_n$ , the lists  $L_1, \dots, L_k$ , and the parameter  $c > 0$ , and do one of the following two things.

- Return lists  $L'_1, L'_2, \dots, L'_k$  satisfying properties (A\*) and (B); or
- Report (correctly) that there is no set of lists  $L'_1, L'_2, \dots, L'_k$  that satisfies both properties (A\*) and (B).

## Problem 3.1

### Definitions:

$P_i$  = The # of doctors required to work on day  $i$

$K$  = The # of doctors

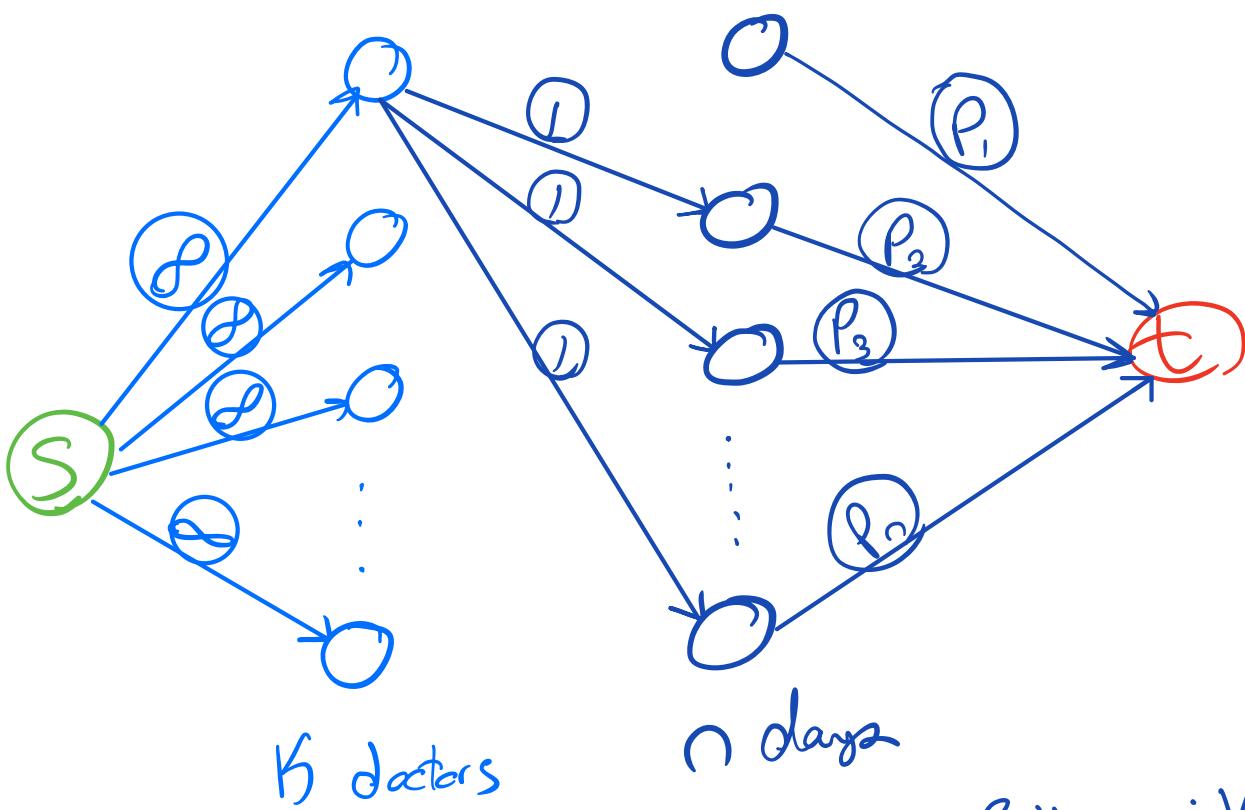
$L_j$  = List of days doctor  $j \leq K$  is willing to work

$L'_j$  = Subset of  $L_j$  if given all  $L_j \rightarrow$  the hospital can obtain

$P_i$  doctors for the next  $n$  days ( $1 \leq i \leq n$ )

Key Idea: Create a <sup>bipartite</sup> graph mapping  $K$ -doctors to the  $n$  days, with edges assigned according to each doctors  $L_j$ . If we can bound edges by capacities according to the problem constraints, we can create a flow-graph and create an assignment via flow.

### Graph Construction:



Define the above graph  $G$  with the following:  $V \rightarrow$  the set of vertices,  $E \rightarrow$  the set of edges.  
 $V = \{S, t\} \cup \{j \mid 1 \leq j \leq K\} \cup \{i \mid 1 \leq i \leq n\}$

$E = \{(S, j) \mid \forall j \in K\text{-doctors}\}$  with a conceptual capacity of  $\infty$ .  
 $\cup \{(j, i) \mid \forall j \in K\text{-doctors} \text{ and } i \in n\text{-days with an edge only created if day } i \text{ is in doctors } L_j\}$   
 with a capacity of 1.  
 $\cup \{(i, t) \mid \forall i \in n\text{-days}\}$  with a capacity of  $P_i$ ,  
 the number of doctors required to work on day  $i$ .

### Algorithm:

Create a bipartite graph  $G$  as defined above in Graph Construction. Compute the max flow using 1 of the 4 max flow algorithms (for selection see Time Complexity and Max Flow algorithms). If the max flow =  $\sum_{i=1}^n P_i$ , then create  $L'_j \forall j \in K$  doctors. Specifically these  $L'_j$  equate to the edges from doctor  $j$  to all days in  $L_j$  with non-negative flow (which here will be 1). If the max flow <  $\sum_{i=1}^n P_i$ , then report that there is no set of lists  $L'_1, L'_2, \dots, L'_K$  that satisfy the problem constraints.

Implementation Intricacy: The maximum flow can be checked by summing the flow out of the source node  $S$ , along its outgoing edges.

## Proof of Correctness:

### CLAIM 1:

$\exists$  an assignment  $L'_1, L'_2, \dots, L'_k$  from the lists  $L_1, L_2, \dots, L_B$ , if the maximum flow through  $G = \sum_{i=1}^k P_i$ . Otherwise no assignment exists.

- There are 2 capacity bottlenecks in  $G$ :
  - 1) Capacity  $(P_i)$  for all  $i \in n$ -days on  $(i, t)$  edges
  - 2) Capacity  $(1)$   $\forall (j, i)$  edges, which were created only if doctor  $j$  is willing to work on day  $i$  (i.e., day  $i$  is in  $L_j$ )
- Given this,  $(j, i)$  edges establish whether or not a doctor can work on day  $i$ , contributing at most 1 unit of flow.
- For a day  $i$ , there needs to be  $P_i$  units of flow  $\sum_{(j,i)} 1$  the incoming flow. This ensures that every  $(j, i)$  edge with flow (note the flow along these edges can only be 0 or 1 based on these edge capacities) equates to 1 doctor working on day  $i$ . Thus all edges with flow must  $= P_i$  for all  $n$  days for there to be a valid assignment  $L'_1, L'_2, \dots, L'_k$ .
- The capacities on  $(i, t)$  edges ensure we cannot overflow a particular  $(i, t)$  edge ensuring that when  $\text{maxflow} = \sum_{i=1}^k P_i$  every  $(i, t)$  edge is saturated with exactly  $P_i$  flow.

• Note: This also effectively constrains the maximum flow in  $G$  to be  $\sum_{i=1}^n p_i$

• When  $\text{maxflow} < \sum_{i=1}^n p_i$ , that equates to the fact that there is at least one day where  $p_i$  doctors cannot work.

### CLAIM 2:

If the maximum flow =  $\sum_{i=1}^n p_i$ , then all  $(j, i)$  edges with flow correspond to a valid assignment. Specifically, for any doctor  $j$ , these

edges to days  $i$ , with flow correspond to  $L'_j$ .

• By CLAIM 1, know that flow can only =  $\sum_{i=1}^n p_i$  when there is an assignment.

• Since flow is bounded to be either 0 or 1 on these edges, by the capacity (assuming  $p_i$  is an integer... hard to have half a doctor work on a day), all edges with flow (i.e., flow = 1 since flow will be pushed along edge  $(i, t)$ ) in discrete units of 1.

• Note: If  $p_i$  is not an integer, assume hospital needs a partial shift. In such a case all individual  $(j, i)$  edges w/  $\text{flow} = 1$  equate to full shifts and all  $(j, i)$  edges with  $\text{flow} > 0 & < 1$  equating to partial shifts (i.e., 0.2 flow equates to a 20% shift)

• This would necessitate obtaining a least common multiple of all capacities to convert all capacities in  $G$  to integers, before scaling the flow back by such a multiple.

- Note that when  $\text{flow} = 0$  in the existence of a valid assignment, this corresponds to a day  $i$  that is saturated w/  $P_i$  doctors, & thus no more doctors need to be assigned to day  $i$ . Thus doctors only work on flow  $(i, i)$  edges.

### Time Complexity and Max Flow Algorithm Selection:

Creation of  $G$  and the individual steps of the max-flow algorithm are linear in the size of the graph. Thus the max-flow algorithm selection is the time complexity bottleneck.

Note in the worst case, every doctor could work every day creating the following size of  $G$ :

$$|V| = K + n + 2$$

$$|E| = K + n + Kn$$

Given the dense-ness of  $|E|$  in this case, omit algorithms non-linear in  $E$  (i.e., Scaling and KED).

Preflow-push gives:  $O((K + n + 2)^3)$  which contains not only  $Kn$ , but also  $Kn^2 + K^2n$

Ford-Fulkerson gives  $O((2K + 2n + Kn + 2) \sum_{i=1}^2 P_i)$ , where  $\sum_{i=1}^2 P_i$ , the max possible flow is upper bounded by  $Kn$  (i.e., the number of doctors \* the number of days; if any  $P_i > K$  can terminate early as know we cannot assign more doctors than we have available.)

Substituting  $K_n$  to  $\sum_{i=1}^n P_i = K_n$  to the Ford-Fulkerson algorithm get  
 time complexity of:

$$O((2K+2n+3n+2) \underbrace{\sum_{i=1}^n P_i}_{Kn}) = O(2K^2n + 2Kn^2 + K^2n^2 + 2Kn)$$

Compare this to Pre-flow push which has  $K^3, n^3, K_n^2, K^2n$   
 terms, but its exponents are bounded by 3. Thus in the worst case  
 where  $K=n$  the highest order term in Ford-Fulkerson is  $n^4$  vs.  
 in pre-flow push it is  $n^3$ .

Thus we select the pre-flow push algorithm here with a  
 time complexity of  $O((K+n+2)^3) \approx O((Kn)^3)$ .

### Problem 3.2:

#### Definitions:

$P_i, K, L_j$  are defined in the same manner as in Problem 3.1  
 above.

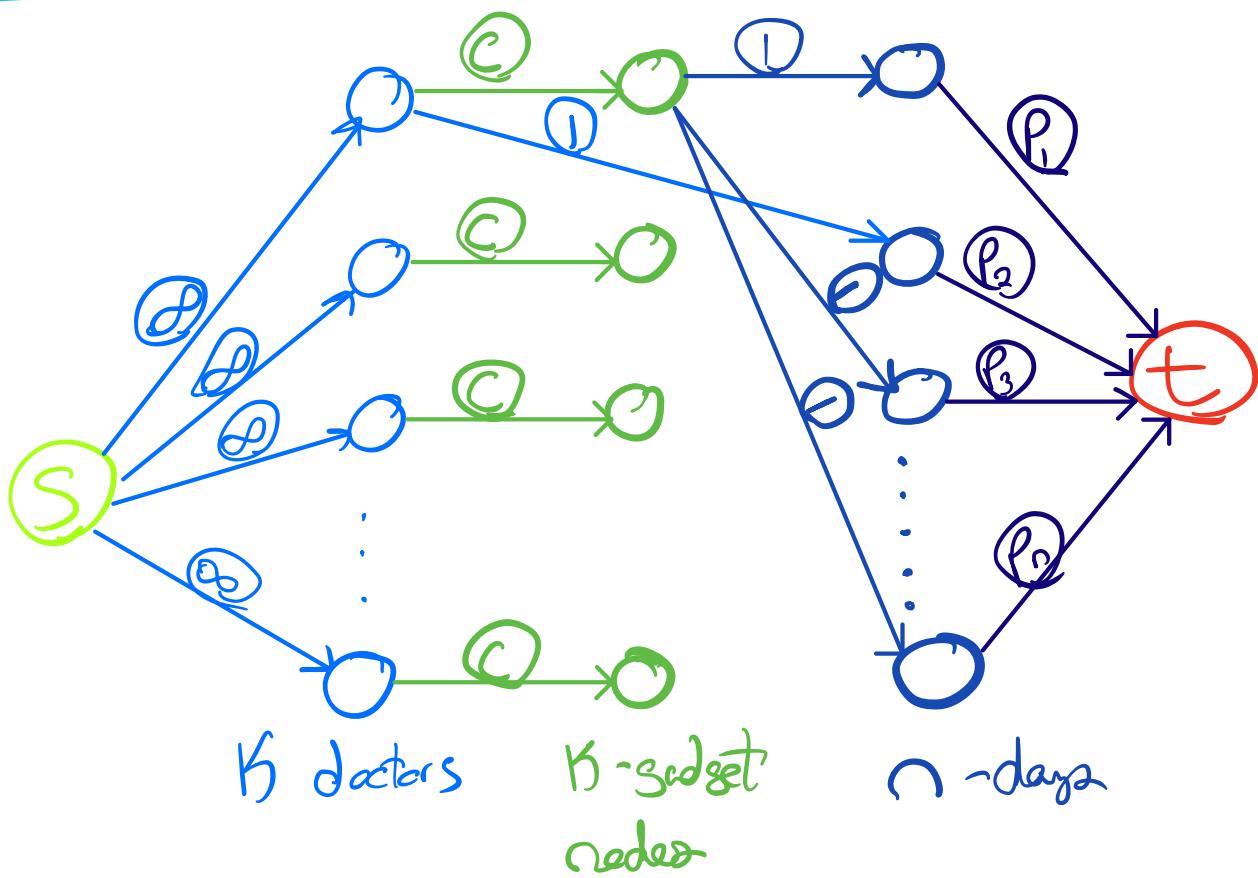
$C$  = flexibility parameter  $> 0$  s.t. given a doctors list  $L_j$ , doctor  
 $j$  can be assigned to work no more than  $C$  days outside of his/her  
 $L_j$ .

$L'_j$  = Subset of  $L_j$  + at most  $C$  days not in  $L_j$  if given all  
 $L_i \rightarrow K$  the hospital can obtain  $P_i$  doctors for the next  $n$  days.

## Key Idea:

Introducing the flexibility parameter  $C$  introduces a new constraint that needs to be accounted for if attempt to build on flow-approach from Problem 3.1 above. To account for this constraint we can introduce another layer of  $K$  "gadget" nodes with an inbound capacity of  $C$ .

## Graph Construction:



Define the above graph  $G$  with the following:  $V \rightarrow$  the set of vertices,  $E \rightarrow$  the set of edges.

$V = \{s, t\} \cup \{j \mid j \in K\text{-doctors}\} \leftarrow \text{doctor nodes}$  $\cup \{g_j \mid j \in K\text{-doctors}\} \leftarrow \text{gadget nodes}$  $\cup \{i \mid i \in n\text{-days}\}$  $E = \{(s, j) \mid \forall j \in K\text{-doctors}\}$  with a conceptual capacity of  $\infty$ . $\cup \{(j, g_j) \text{ for all } K\text{-doctors, s.t. there is only one edge between doctor } j \text{ and gadget } g_j \text{ (the indegree of } g_j \text{ will only be 1 with the lone edge coming from doctor } j)\}$  with a capacity of  $C$ . $\cup \{(j, i) \mid \forall j \in K\text{-doctors and } i \in n\text{-days with an edge only created if day } i \text{ is in doctors } L_j\}$  with a capacity of 1. $\cup \{(g_j, i) \mid \forall i\text{-days for doctor } j \text{ for which no } (j, i) \text{ edge exists}\}$  with a capacity of 1. $\cup \{(i, t) \mid \forall i \in n\text{-days}\}$  with a capacity of  $P_i$ , the number of doctors required to work on day  $i$ .

## Algorithm:

Create a bipartite graph  $G$  as defined above in Graph Construction. Compute the max flow using 1 of the 4 max flow algorithms (for selection see Time Complexity and Max Flow algorithms). If the max flow =  $\sum_{i=1}^n p_i$ , then Algorithm Selection below). If the max flow <  $\sum_{i=1}^n p_i$ , then create  $L'_j \forall j \in K$  doctors. Specifically these  $L'_j$  equate to the edges from doctor  $j$  to all days in  $L_j$  with non-negative flow (which here will be 1) plus all the edges from doctor  $j$ 's corresponding gadget node to all  $i$  days w/ which it is connected with non-zero flow. If the maximum flow <  $\sum_{i=1}^n p_i$ , then report that there is no set of lists  $L'_1, L'_2, \dots, L'_K$  that satisfy the problem constraints A\* and B.

Implementation Intricacy: The maximum flow can be checked by summing the flow out of the source node  $S$ , along its outgoing edges.

## Proof of Correctness:

### CLAIM 1:

$\exists$  an assignment  $L'_1, L'_2, \dots, L'_K$  from the lists  $L_1, L_2, \dots, L_K$ , if the maximum flow through  $G = \sum_{i=1}^n p_i$ . Otherwise no assignment exists.

- There are 3 capacity bottlenecks in  $G$ :

- 1) Capacity  $P_i$  for all  $i \in n$ -days on  $(i, t)$  edges
- 2) Capacity  $1$   $\forall (j, i)$  edges, which were created only if doctor  $j$  is willing to work on day  $i$  (i.e., day  $i$  is in  $L_j$ ). Similarly there are capacities of  $1 \forall (g_j, i)$  edges, which were created only on days for which doctor  $j$  was unwilling to work.
- 3) Capacity  $C$   $\forall (j, g_j)$  edges, the flexibility parameter for doctors which cannot be exceeded.

- Know that Capacity bottleneck 1) above ensures a prevention of flow oversaturation. Specifically the capacity of  $P_i$  here ensures at most  $P_i$  units can flow along a particular  $i \rightarrow t$  path. Thus if flow equals  $\sum_{i \in n} P_i$  then for each day  $i$ , the flow along its corresponding edge in  $G$  is exactly  $P_i$  corresponding to the problem constraints.
- The capacities on  $(i, t)$  edges ensure we cannot overflow a particular  $(i, t)$  edge ensuring that when  $\text{maxflow} = \sum_{i=1}^n P_i$  every  $(i, t)$  edge is saturated with exactly  $P_i$  flow.

- Similarly note that capacity bottleneck 2) above ensures a 1-1 doctor to day amount of flow. From any gadget or doctor node the flow is upper bounded by 1, and thus every edge here with flow contributes 1 unit to  $\rho_i$ , which equates to the problem, every doctor counts once to the sum  $\rho_i$  (assuming no fractional requirements; if desired this can be accounted for → See note in Problem 3.) above).
- Capacity Constraint 3) is needed to satisfy the flexibility requirement of the problem & is justified in the reasoning for CLAIM 2 below. However note that all paths that flow through these gadget nodes can have at most one-unit of flow as all outgoing edges from  $G_3$  have a capacity of 1 as stated in the above point.
- When  $\text{maxflow} < \sum_{i=1}^n \rho_i$ , that equates to the fact that there is at least one day where  $\rho_i$  doctors cannot work. This occurs when all incoming edges to  $i$  are capacity saturated by  $C$ . Specifically since there are  $K$  edges incident on all  $n$  days, all non-gadget edges will be assigned flow to day  $i$ , & thus this necessarily means that capacity bottleneck 3) caps the flow to day  $i$  as all  $g_j \rightarrow i$  edges are possible, but no flow runs through any of them, → no assignment exists that satisfies the problem constraints A\* & B.

## CLAIM 2:

If the maximum flow =  $\sum_{i=1}^n p_i$ , then all  $(j, i)$  and  $(g_j, i)$  edges

with flow correspond to a valid assignment of doctor  $j$  to  $i$  days,  
Corresponding to their  $L'_j$ .

- For a valid assignment to exist all  $i \in n$ -days, must have  $p_i$  doctors assigned on day  $i$ 
  - By CLAIM 1, know that flow can only =  $\sum_{i=1}^n p_i$  when there is an assignment.
- For a valid assignment to exist, doctors should only be assigned to days in their corresponding list,  $L'_j$ , and at most  $C$ -days outside that list.
  - By the capacity constraint 2) know that the flow through edges outlined in capacity constraint 2) account to an assignment as flow is upper-bounded by 1 along these paths (for more detailed reasoning see above).
  - While this is valid, need to ensure an excess of flow cannot run through the gadget nodes. If this were to occur this would violate the flexibility constraint imposed by the problem.
  - Capacity Constraint 3) ensures that at most  $C$  units of flow can run through a gadget node. This ensures the flexibility constraint is adhered to and thus by conservation of flow the exiting flow cannot exceed the incoming flow which is bounded by  $C$ .
- Thus a valid assignment can be obtained as outlined by this Claim.

## Time Complexity and Max Flow Algorithm Selection:

Creation of  $G$  and the individual steps of the max flow algorithm are linear in the size of the graph. Thus the max flow algorithm selection is the time complexity bottleneck.

Given the creation of the graph observe the following:

$$|V| = 2k + n + 2 \leftarrow \text{gadget} + \text{doctors} + \text{days} + \text{source} + \text{sink}$$

$$|E| = k + k + k_n + ? = k(2+n) + n = 2k + n + kn$$

$$|E| = \underbrace{k}_{\text{sources to doctors}} + \underbrace{k}_{\text{doctors to gadget}} + \underbrace{k_n}_{\text{gadgets to days}} + \underbrace{?}_{\text{doctors to days}} \rightarrow \text{days to sink}$$

Note this is within constant bounds of the number of vertices + edges in problem 3.1 above with the following:

$$|V| = k + n + 2$$

$$|E| = k + n + kn$$

Thus by the same reasoning outlined in the justification stated in Problem 3.1, select pre-flow push here yielding a time complexity of:  $\mathcal{O}((2k + n + 2)^3) \approx \mathcal{O}((k + n)^3)$

#### Problem 4: Cellular network

Consider the problem of selecting nodes for a cellular network. Any number of nodes can be chosen from a finite set of potential locations. We know the cost  $c_i \geq 0$  of establishing site  $i$ . If sites  $i$  and  $j$  are selected as nodes, then we derive the benefit  $b_{ij}$ , which is the revenue generated by the traffic between the two nodes. Both the benefits and costs are non-negative integers. Find an efficient algorithm to determine the subset of sites as the nodes for the cellular network such that the sum of the benefits provided by the edges between the selected nodes less the selected node costs is as large as possible.

Design an efficient polynomial-time algorithm.

Provide a high-level description of your algorithm, prove its correctness, and analyze its time complexity.

Goal:

Find the subset of nodes to maximize profits. Specifically

$$\max_{S \subseteq \text{Cellular nodes}} \left( \sum_{i,j \in S} b_{i,j} - \sum_{i \in S} c_i \right)$$

Definitions:

$c_i$  = the cost of establishing site  $i$  with an integer cost  $\geq 0$ .

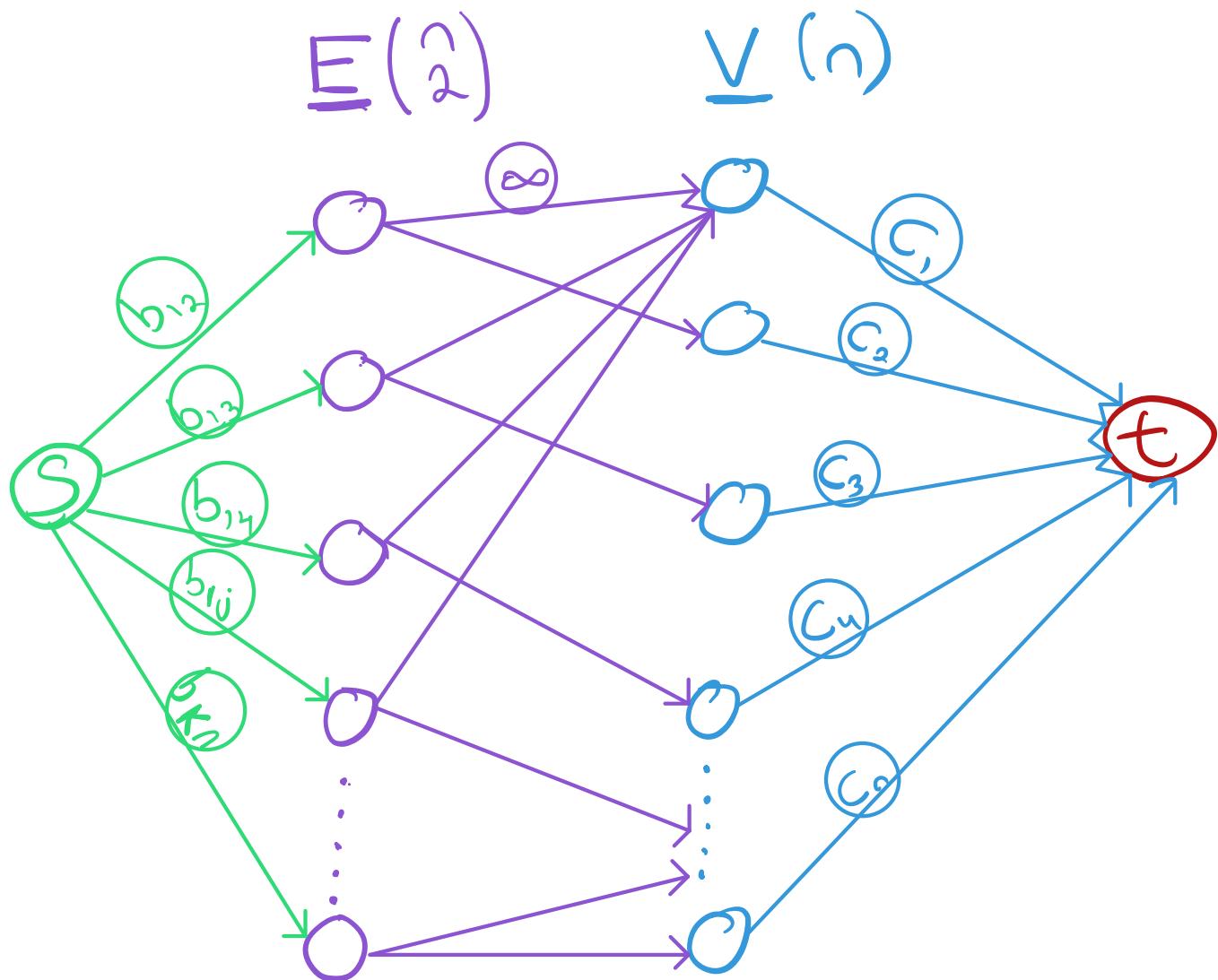
$b_{ij}$  = benefit derived by the traffic between nodes  $i$  &  $j$   
 This benefit is a non-negative integer

Profit = Sum of all benefits across all pairs of established nodes - the sum of the costs of establishing all the selected nodes

Key Idea:

Cellular nodes create an undirected graph, which we can convert to a bipartite graph by conceptualizing undirected edges as vertices. Moreover, by assigning benefits and costs correctly in the context of a flow graph & running a max flow algorithm, the subset of nodes to select that maximizes profits equate to the nodes reachable from  $S$  in  $R_G \Rightarrow \min \text{ cut}$ .

## Graph Construction:



Define the above graph  $G$  with the following:  $V \rightarrow$  the set of vertices,  $E \rightarrow$  the set of edges in the bipartite graph

$V = \{S, t\} \cup \{ \text{all (finite) potential locations for cellular node site} \}$   $\leftarrow$  define the number of sites as  $n$   
 $\leftarrow$  this is  $V$  in the above figure

$\cup \{ \text{all } (i,j) \text{ edges that exist in the undirected graph} \}$

- Note that there are  $\binom{n}{2}$  edges in the undirected graph thus this term contributes  $\binom{n}{2}$  vertices
- This is  $E$  in the above figure

$E = \{(s, e_{ij}) \text{ for all edges that exist in the undirected graph}\} \text{ with a capacity of } b_{ij} \rightarrow \text{the benefit derived from sites } i \text{ and } j$

$\cup \{(e_{ij}, V_i), (e_{ij}, V_j)\}$  with a conceptual capacity of  $\infty$ .

- This states that each  $e_{ij}$  only connects to its corresponding vertices  $V_i$  and  $V_j$ . Note  $i \neq j$ .

$\cup \{(V_i, t) \in \text{potential locations}\}$  with a capacity of  $c_i$ , the cost of establishing site  $i$ .

### Algorithm:

Create a bipartite graph  $G$  as defined above in Graph Construction above. Run one of the four max flow algorithms (for selection see Time Complexity and Max Flow Algorithm Selection section below). The maximum profit that can be obtained equates to the min cut (which equates to the max flow by duality) on  $G$ . The subset of sites as nodes for the cellular network equate to the subset of nodes reachable from  $S$  in the Residual Graph  $R_G$ , at the time max flows reached in  $G$  (i.e., when there are no more Augmenting Paths in  $R_G$ ). In  $R_G$ , these sites can be determined by running BFS from  $S$  and returning all reachable sites (i.e.,  $V_i$ ).

## Time Complexity and Max Flow Algorithm Selection:

Creation of  $G$  and the individual steps of the max flow algorithm are linear in the size of  $G$ . Running BFS on  $R_G$  is also linear in the size of  $G$ . Thus the time complexity bottleneck is the max flow algorithm selected.

Given  $G$  as defined in Graph Construction observe the following size of  $G$ :

$$|V| = \binom{n}{2} + n + 2 = \frac{n(n-1)}{2} + n + 2 \approx n^2$$

$$|E| = \binom{n}{2} + 2n + n = \frac{n(n-1)}{2} + 3n \approx n^2$$

Using the approximation of both  $|V|$  and  $|E|$  as their highest order terms  $n^2$  we observe the following time complexities for each of the max flow algorithms

Ford-Fulkerson:  $2n^2 \cdot \text{max flow (or min cut)}$

• Without a clear upper bound on max flow here we can use the less tight bound  $C_{\text{total}} = \sum_{e \in E} C_e$  for all non-infinite capacity edges.

Scaling:  $(n^2 + n^2)^2 \log C_{\text{max}} = 4n^4 \log C_{\text{max}}$

•  $C_{\text{max}}$  is the maximum non-infinite capacity in the graph  $G$  here

$$\text{Karp-Edmonds-Dinitz: } (n^2 + n^2)^2 \cdot n^2 = 4n^6$$

$$\text{Preflow push: } (n^2)^3 = n^6$$

Thus the choice of algorithm here depends on  $C_{\max}$  and  $C_{\text{total}}$ . Note we can discard KDE as it has a worse time complexity than Preflow push by a constant factor. For the 3 remaining algorithms, observe they all share an  $n^2$  term, so need to look at the remaining defining characteristics. Specifically the minimum among  $C_{\text{total}}$ ,  $4n^2 \log C_{\max}$ , and  $n^4$  determine which method is selected. If  $C_{\text{total}}$  is the minimum, select Ford-Fulkerson. If  $4n^2 \log C_{\max}$  is the minimum, select Scaling. If  $n^4$  is the minimum select Preflow Push.

Given the dependency on the relationship between the number of sites, costs, and benefits, default to the lowest order polynomial algorithm: Ford-Fulkerson. This gives a time complexity of  $\mathcal{O}(2n^2 \cdot C_{\text{total}}) \rightarrow \mathcal{O}(n^2 C_{\text{total}})$ .

## Proof of Correctness:

### CLAIM 1:

Given  $G$ , the subset of sites extracted from the minimum cut in  $R_G$ , equates to the subset of sites that maximize profit.

Let  $A$  be the set of vertices in  $R_G$  s.t. they can be reached from the source node  $S$ .

- Let  $B$  be the complement of  $A$ , that is  $B = V - A$
- Let  $B$  be the complement of  $A$ , that is  $B = V - A$
- Know that at the completion of max flow algorithm  $CAP(A, B) = \text{maxflow}$
- $CAP(A, B) = \sum_{ij \notin A} b_{ij} + \sum_{i \in A} c_i$  ← note omission of  $\infty$  capacity edges. See CLAIM 2 below.

Define  $b = \sum_{ij} b_{ij}$  ← this equates to the cut capacity of just  $A = S$ , and is a constant value.

Rearrange  $CAP(A, B)$  by adding and subtracting  $\sum_{ij \in A} b_{ij}$ :

$$CAP(A, B) = \underbrace{\sum_{ij \notin A} b_{ij} + \sum_{ij \in A} b_{ij}}_b - \sum_{ij \in A} b_{ij} + \sum_{i \in A} c_i$$

$$\begin{aligned}
 \text{CAP}(A, B) &= K - \sum_{i,j \in A} b_{ij} + \sum_{i \in A} c_i \\
 &= K - \underbrace{\left[ \sum_{i,j \in A} b_{ij} - \sum_{i \in A} c_i \right]}_{\text{Profit}}
 \end{aligned}$$

Thus  $\text{CAP}(A, B)$  is minimized when profit is maximized, and thus the min cut and all nodes in  $A$  correspond to the subset of sites that maximize the profit.

## CLAIM 2:

For any min cut that occurs, there must necessarily be no cut edges with  $\infty$  capacity.

Cut edges with  $\infty$  capacity from

- In this case a cut will either occur across edges from  $S$  or inband to  $t$ .
- This can be shown by observing that there always exists a finite capacity cut in  $G$  by selecting from these edges + always omitting  $\infty$  capacity edges.
- This Claim has important implications for the cellular site selection. Specifically that all projects selected in  $A$  have at least one path from  $S$  in  $R_G$  + thus their cost is the flow bottleneck → Cost must necessarily  $\geq$  benefit.

### CLAIM 3:

Given the cellular node site selection project, it can be reduced to a flow graph  $G$ , that satisfies the constraints of the problem.

- The capacity bottleneck is equating to pairwise node benefits and costs, coupled with unbounded capacity edges between  $E$  and  $V$  yield the maximum profit after running the selected max flow.
- The detailed reasoning for this is shown above in CLAIM 1.