



SIGMOD Programming Contest 2014

Large Social Network Analysis

Team: blxlrsmb
Tsinghua University



Task Overview

Implement a **social network analysis** system.

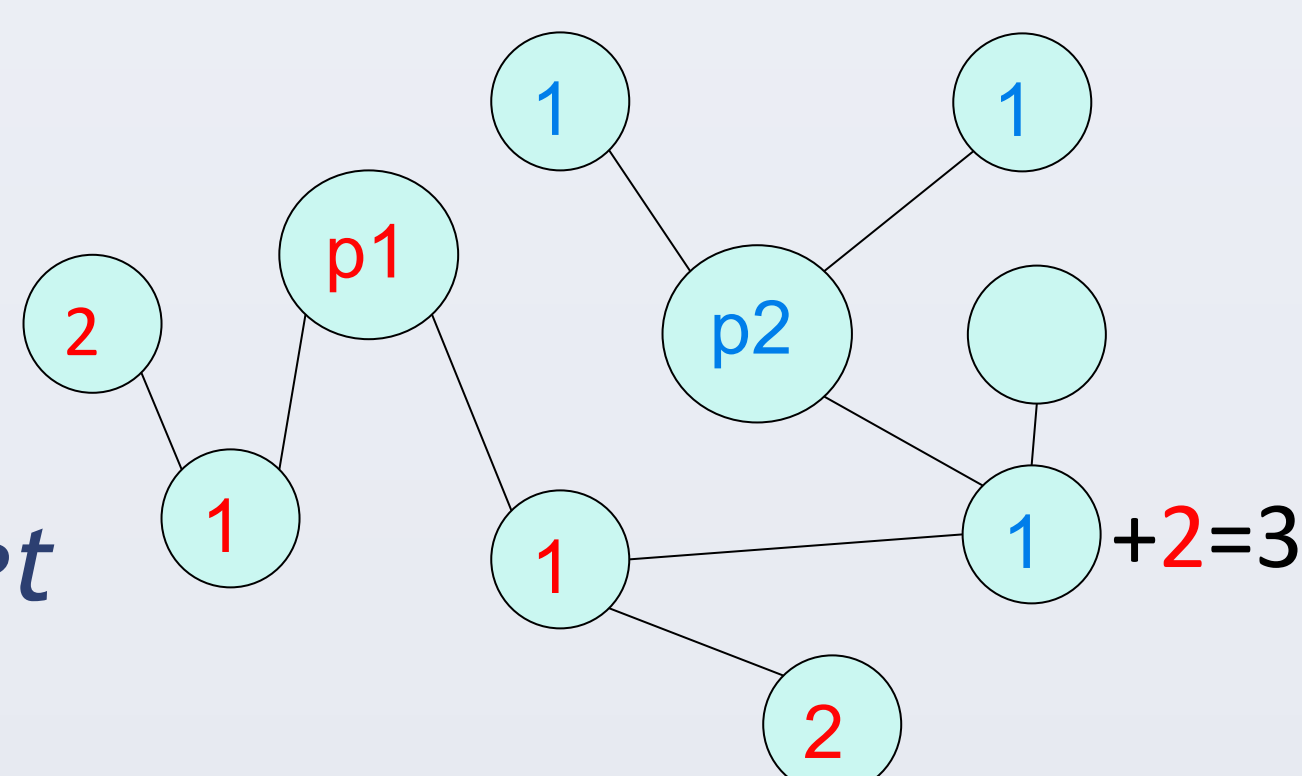
Given a huge social network as an undirected weighted graph $G(V, E)$ with the following attributes on each vertex v : a numeric attribute v_d , a set of tags v_t , we should implement an analysis system which reads dataset & 4 types of queries and produce the results.

Query 1

Given x and vertex p_1, p_2 , calculate the length of shortest path from p_1 to p_2 s.t. each edge e on the path has weight $e_w \geq x$.

Solution:

Execute BFS from both p_1 and p_2 . Use a bit-array maintaining all the visited vertices, to see if they have met in the middle.



Query 2

Given k and d , find top- k tags with largest range. Here the range of a tag T in graph $G(V, E)$ is defined as the size of the largest subgraph $G(V', E')$ s.t. $\forall v \in V', v_d > d \text{ \& } T \in v_t$.

Solution:

First **sort all the queries** and vertices by c in descending order. Build an empty graph G_0 and incrementally insert sorted vertices and their corresponding edges to G_0 .

During the insertion, we maintain the top- k largest connected components for each tag using **Union-Find Set**. A query q would finish as soon as all vertices v with $v_d > q_d$ are inserted.

Query 3

Given k, h, p , where p is used to define a subset V' of V . Find the top- k pairs of vertex (u, v) ordered by $\|u_t \cap v_t\|$, also satisfying: $u, v \in V \text{ \& } dist(u, v) \leq h$.

Solution:

First we build V' according to p . Then we build an inverted list for every tag T : $L(T)$ = a list of vertex $v \in V'$ where $T \in v_t$.

This way, we can quickly find vertices in V' sharing tags with a given v by counting # of **occurrences in inverted lists**.

Other Tricks

- Use `mmap()` to read data files faster.
- Implement a thread pool to utilize multi-core processors.
- Google's `tcmalloc` is faster to allocate memory.
- Another way of BFS in estimating lower bound in query 4 is to use a bitset $B_l(v)$ of length $\|V\|$ setting all the vertices u where $dist(u, v) = l$. Since $B_{l+1}(v) = \neg B_l(v) \wedge \bigvee_{(u,v) \in E} B_l(u)$, the search can be implemented using **SIMD instructions**.

Query 4

Given k, t , where t is used to define a subgraph G' of G , find the top- k centralized vertex in G' . Here the centrality of a vertex v in a graph $G(V, E)$ is defined as :

$$C(v) = \frac{r(v)^2}{(\|V\| - 1) \sum_{u \in V} dist(u, v)}$$

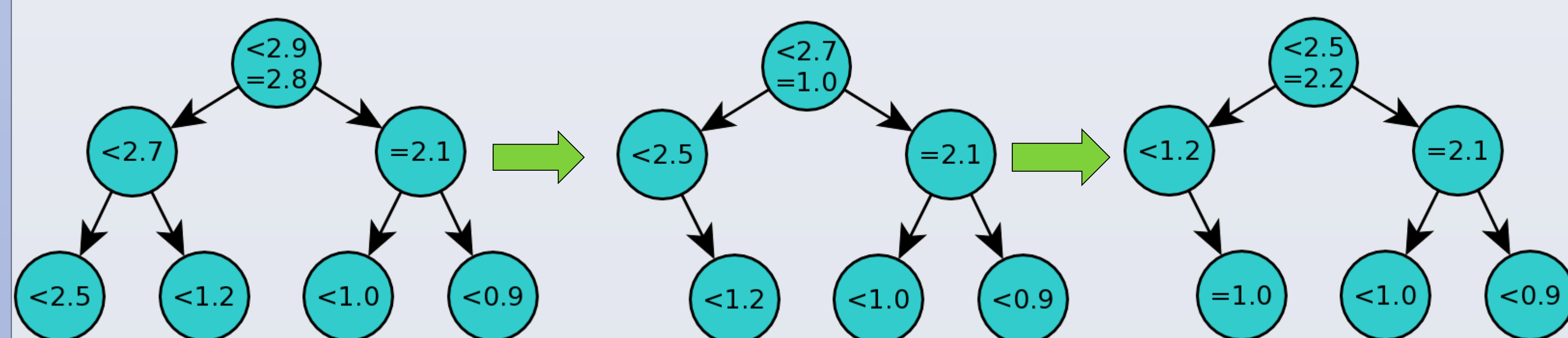
where $r(v)$ is # of vertices reachable from v (exclusive).

Solution:

The subgraph can be efficiently built. The difficulty lies in the calculation of $C(v)$, since $S(v) = \sum_{u \in V} dist(u, v)$ has $O(\|V\|^2)$ time complexity for each v . The solution shall make full use of 'top- k ' in pruning as well as the properties of social network to be efficient.

General framework to obtain top- k centrality:

1. Estimate an **upper bound** of $C(v)$ (from lower bound of $S(v)$) for each v and maintain them with a **max-heap**.



2. Iteratively pop the top element from the heap and calculate its actual value. If the actual value is still the largest in heap, then it's certainly the largest among all actual values of remaining element. Otherwise we insert this actual value back to the heap.

Estimation of lower bound of $S(v)$:

$$S(v) \geq \sum_{dist(u,v) \leq l} dist(u, v) + \sum_{dist(u,v) > l} (l + 1)$$

The first term is calculated by a BFS from v with search depth limited to l . The second term is obtained by counting.

To choose a proper l , we compare our current estimation of each $S(v)$ to an **approximation** of $S(v)$ and increase l when the mean-square error is large. Since **social networks have small diameter (!)**, a small l (3 or 4) is sufficient to guarantee a good lower bound estimation of $S(v)$.

Approximation of $S(v)$:

Random sample a subset of V called V' , and run a thorough BFS on each of them. Then for each $v \in V, S(v) \approx \sum_{u \in V'} dist(u, v) \frac{\|V\|}{\|V'\|}$. This is a highly accurate approximation, which gave an average 1~2% error sampled with 0.1% of V in our experiments.

This approximation can be applied in pruning: by calculating actual centrality of the vertices with top- $(3k)$ approximated centrality, we get $3k$ actual centrality that is very likely to contain the actual top- k centralities. Then the k th centrality among the $3k$ results can give a strong pruning.

Team Members:

Yuxin Wu
ppwwyyxx@gmail.com
Jiawen Liang
taobingxue001@126.com

Xinyu Zhou
zxytim@gmail.com
Junbang Liang
williamm2006@126.com

Wenbo Tao
thierryhenrytracy@163.com
Han Zhao
nikifor383@gmail.com