

DISTRIBUTED MINIMUM SPANNING TREES

CME 323 Project

Swaroop Indra Ramaswamy & Rohit Patki

May 27, 2015

Stanford University

PROBLEM STATEMENT & ASSUMPTIONS

Problem Statement

Given an undirected, weighted simple graph we attempt to find a Minimum Spanning Tree (MST) or Minimum Spanning Forest (MSF).

Assumptions

- The number of edges in the graph (m) is much larger than the number of vertices of the graph (n)
- The edges of the graph do not fit in the memory of a single machine
- The vertices of the graph do fit in the memory of a single machine

APPLICATIONS

- Single linkage clustering
- Network Design
- Image Segmentation
- Taxonomy
- Broadcasting in computer networks
- Important primitive in many graph algorithms

Classical Algorithms

- Kruskal's Algorithm - $O(m \log n)$
- Prim's Algorithm - $O(m \log n)$

Faster Algorithms

- Karger, Klein and Tarjan (1995) - Randomized $O(m)$
- Bernard Chazelle (2000) - $O(m\alpha(m, n))$
- Fredman and Willard (1994) - $O(m + n)$ for integer weights

EDGE PARTITIONING

if $|E| < \eta$ **then**

 Compute $T^* = MST(V, E)$

 Return T^*

end if

$k = \Theta\left(\frac{|E|}{\eta}\right)$

Partition E into E_1, E_2, \dots, E_k where $|E_i| < \eta$ using a universal hash function

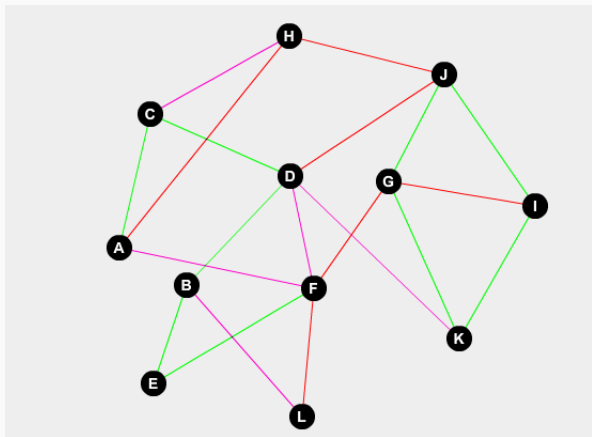
Compute $T_i^* = MST(G(V, E_i))$ in parallel

Return $MST(G(V, \cup_i T_i^*))$

Algorithm 1: Edge partitioning

EDGE PARTITIONING

Figure: Edge Partition



EDGE PARTITIONING : ANALYSIS

- Memory of each machine, $\eta = O(n^{1+\epsilon})$
- Number of edges, $m = O(n^{1+c})$
- Number of machines, $k = O(n^{c-\epsilon})$

Processing time

It can be shown that the algorithm takes $\lceil \frac{c}{\epsilon} \rceil$ iterations in expectation.

Processing time per iteration

$$\underbrace{O\left(\frac{m}{k} \log n\right)}_{\text{Kruskal's on each machine}} + \underbrace{O\left(\frac{m}{k}\right)}_{\text{random partitioning of edges}}$$

Communication Cost

- One all-to-all communication (shuffle) in each iteration
- Cost per iteration = $m, \frac{m}{n^\epsilon}, \frac{m}{n^{2\epsilon}}$

Total communication cost

$$\frac{n(n^\epsilon - 1)}{1 - n^{-\epsilon}}$$

VERTEX PARTITIONING

Partition V into V_1, V_2, \dots, V_k with $V_i \cap V_j = \emptyset$ using a universal hash function

Denote the edges induced by V_i and V_j by $E_{i,j}$

Denote the induced subgraph by $G(V_i \cup V_j, E_{i,j})$

Compute $T_{i,j}^* = \text{MST}(G(V_i \cup V_j, E_{i,j}))$ in parallel

Return $\text{MST}(G(V, \cup_{i,j} T_{i,j}^*))$

Algorithm 2: Vertex partitioning

VERTEX PARTITIONING : ANALYSIS

- Number of edges, $m = n^{1+c}$
- Number of partitions, $k = n^{\frac{\epsilon}{2}}$
- Number of edges in each mapper = $n^{1+\frac{\epsilon}{2}}$, in expectation

Processing Time

Total processing time

$$\underbrace{O\left(\frac{m}{k} \log \frac{n}{k}\right)}_{\text{Kruskal's on each machine}} + \underbrace{O\left(n^{1+\frac{\epsilon}{2}} \log n\right)}_{\text{Final Kruskal's}}$$

VERTEX PARTITIONING : ANALYSIS

Communication Cost

- One one-to-all communication (broadcast) = $O(nk^2)$
- One all-to-all (groupByKey) = $O(m)$
- One all-to-one to compute the final MST = $O(n^{1+\frac{\epsilon}{2}})$

Total communication cost

$$O(nk^2) + O(m) + O(n^{1+\frac{\epsilon}{2}})$$

Communication Time

Total communication time

$$O(n \log k) + O(m) + O(n^{1-\frac{\epsilon}{2}})$$

PARALLEL PRIM'S ALGORITHM

A = DISJOINTSET()

for i in V **do**

 A.MAKE-SET(i)

end for

Broadcast A

Find the minimum edge leaving each disjoint set using a reduce operation, denote this by the list of edges, \hat{E}

while $|\hat{E}| > 0$ **do**

for e in \hat{E} **do**

 A.UNION(u, v)

end for

Broadcast A

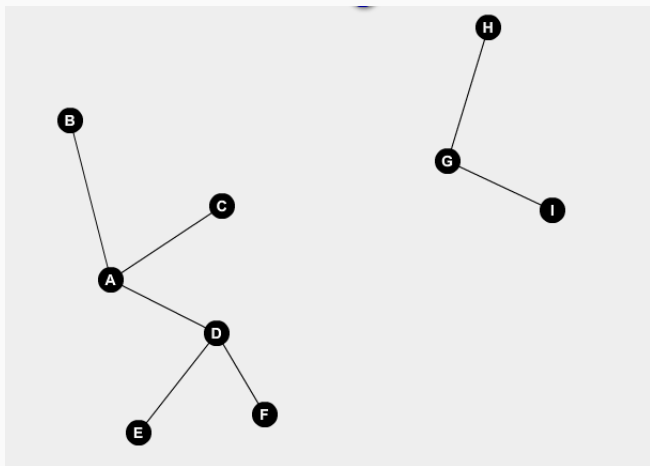
\hat{E} = Minimum edges leaving the disjoint sets

end while

Algorithm 3: Parallel Prim's Algorithm

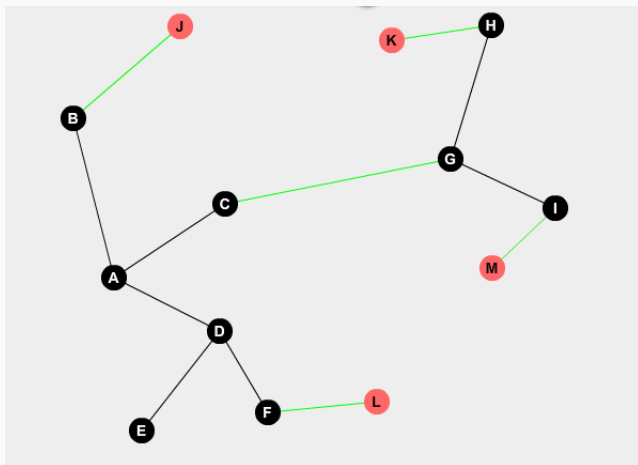
PARALLEL PRIM'S

Figure: Connected components



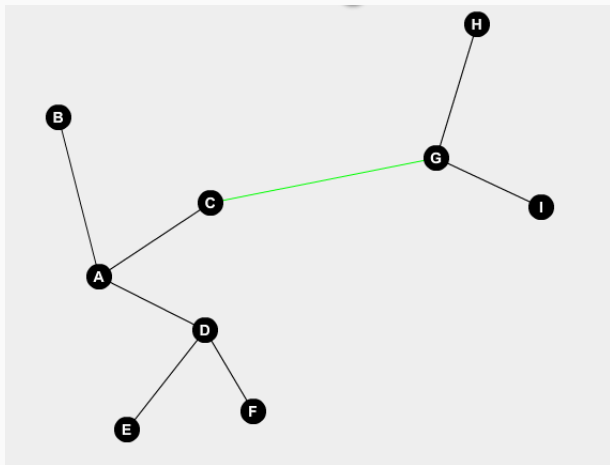
PARALLEL PRIM'S

Figure: Potential new edges



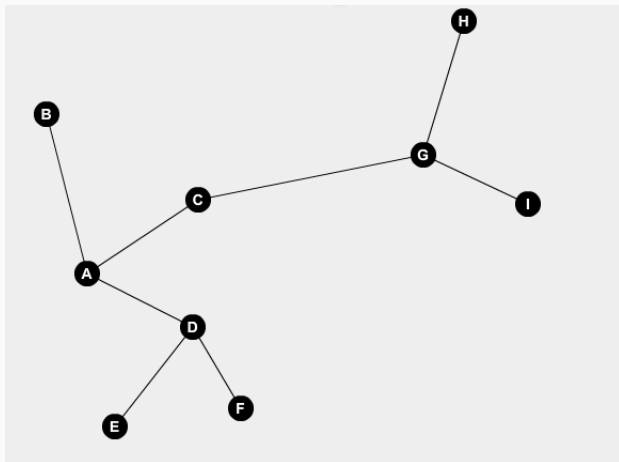
PARALLEL PRIM'S

Figure: Final new edge after reduce



PARALLEL PRIM'S

Figure: New connected component



PARALLEL PRIM'S ALGORITHM : ANALYSIS

After each iteration of the while loop, the number of edges left to find reduces by at least $\frac{1}{2}$. Therefore, at most $\log n$ iterations are required.

Processing Time

Total Processing Time,

$$\underbrace{\log n}_{\text{iterations}} \times \underbrace{O\left(\frac{m}{k}\right)}_{\text{per iteration}} + \underbrace{O(m)}_{\text{Total cost of all the reduces}}$$

PARALLEL PRIM'S ALGORITHM : ANALYSIS

Communication Cost

- One one-to-all broadcast of the disjoint set data-structure, per iteration = $O(nk)$
- One reduce to find minimum edges, per iteration = $O\left(\frac{n}{2^i}\right)$

Total communication cost,

$$O(nk \log n) + O(n)$$

Communication Time

Total communication time,

$$O(n \log k \log n) + O(n \log n)$$

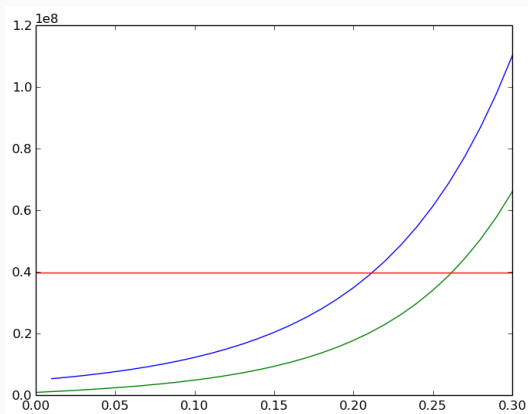
THEORETICAL COMPARISON

	Edge Partitioning	Vertex Partitioning	Parallel Prim's
Proc. Time	$O\left(\frac{m}{\epsilon k} \log n\right)$	$O\left(\left(\frac{m}{n^{\frac{\epsilon}{2}}} + n^{1+\frac{\epsilon}{2}}\right) \log n\right)$	$O\left(\frac{m}{k} \log n + m\right)$
Comm. Time	$O\left(m \frac{1-n^{-c}}{1-n^{-\epsilon}}\right)$	$O(m + cn \log n)$	$O(n \log k \log n)$

- Number of edges $m = n^{1+c}$
- Memory per machine = $n^{1+\epsilon}$

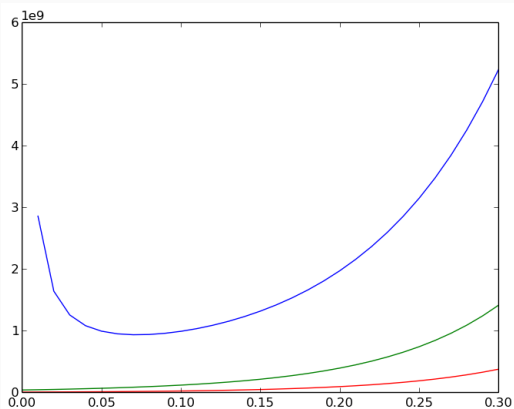
PLOTS: COMMUNICATION TIME

Figure: Communication Time vs vs c for $n = 1,000,000$



PLOTS: PROCESSING TIME

Figure: Processing Time vs c for $n = 1,000,000$



EXPERIMENTAL COMPARISON

	No.of Vertices	No.of Edges	Edge Partitioning	Vertex Partitioning	Parallel Prim's
1	281903	2312497	791 s	316 s	92 s
2	875713	5105039	7384 s	3733 s	229 s
3	685230	7600595	3670 s	1569 s	313 s
4	1696415	11095298	> 7200 s	> 3600 s	335 s
5	1088092	1541898	> 7200 s	> 3600 s	394 s

1. Stanford web graph
2. Google web graph
3. BerkStan graph
4. as-skitter
5. Road-net PA

- 3 algorithms for distributed MST : **Vertex Partitioning**, **Edge Partitioning** and **Parallel Prim's**
- Communication time for Parallel Prim's is independent of the number of edges
- For sparse graphs, if large number of machines are available, use **Vertex Partitioning**
- In general, **Parallel Prim's** has a better processing time than the other two algorithms
- As the density of the graph increases, **Parallel Prim's** wins out