Minimal Spanning Tree

CS60002: Distributed Systems

Pallab Dasgupta Professor, Dept. of Computer Sc. & Engg., Indian Institute of Technology Kharagpur



Leader Election versus Spanning Tree

- Let C_E be message complexity of the election problem and C_T be the message complexity of finding a spanning tree
- Given a spanning tree, we can run election algorithm on tree to find a leader in O(N) time

-Thus: $C_E \leq C_T + O(N)$

 Given a leader, we can run the echo algorithm to find a spanning tree with 2|E| messages

-Thus: $C_T \le C_E + 2|E|$

- Any comparison election algorithm for arbitrary networks has a (worst case and average case) message complexity of at least Ω(|E| + NlogN)
- Therefore the two problems are of the same order of magnitude
 INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

Minimal Spanning Tree

- Let G = (V, E) be a weighted graph, where $\omega(e)$ denotes the weight of edge e.
 - The weight of a spanning tree T of G equals the sum of the weights of the N 1 edges contained in T
 - T is called a *minimal spanning tree* if no spanning tree has a smaller weight than T.

If all edge weights are different, there is only one MST

The notion of a *fragment*

- A fragment is a subtree of a MST
- If F is a fragment and e is the least-weight outgoing edge of F, then F U {e} is a fragment
- Prim's Algorithm:
 - Start with a single fragment and enlarges it in each step with the lowest-weight outgoing edge of the current fragment
- Kruskal's Algorithm:

 Starts with a collection of single-node fragments and merges fragments by adding the lowest-weight outgoing edge of some fragment

- Distributed algorithm based on Kruskal's algorithm
- Assumptions:
 - Each edge e has a unique edge weight $\omega(e)$
 - All nodes though initially asleep awaken before they start the execution of the algorithm. When a process is woken up by a message, it first executes the local initialization procedure, then processes the message

- Algorithm Outline:
 - 1) Start with each node as a one-node fragment
 - 2) The nodes in a fragment cooperate to find the lowest-weight outgoing edge of the fragment
 - 3) When the lowest-weight outgoing edge of a fragment is known, the fragment will be combined with another fragment by adding the outgoing edge, in cooperation with the other fragment
 - 4) The algorithm terminates when only one fragment remains

- Notations and Definitions:
 - 1) Fragment name. To determine whether an edge is an outgoing edge, we need to give each fragment a name.
 - 2) Fragment levels. Each fragment is assigned a level, which is initially 0 for an initial one-node fragment.
 - 3) Combining large and small level fragments. The smaller level fragment combines into the larger level fragment by adopting the fragment name and level of the larger level fragment. Fragments of the same level combine to form a fragment of a level which is one higher than the two fragments. The new name is the weight of the combining edge, which is called the *core edge* of the new fragment.

- <u>Summary of combining strategy</u>: A fragment F with name FN and level L is denoted as F = (FN, L); let e_F denote the lowest-weight outgoing edge of F.
 - Rule A. If e_F leads to a fragment F' = (FN', L') with L < L', F combined into F', after which the new fragment has name FN' and level L'. These new values are sent to all processes in F
 - Rule B. If e_F leads to a fragment F' = (FN', L') with L = L' and $e_{F'} = e_F$, the two fragments combine into a new fragment with level L+1 and name $\omega(e_F)$. These new values are sent to all processes in F and F'.
 - Rule C. In all other cases fragment F must wait until rule A or B applies

- Node and link status:
 - $statch_p[q]$: Node p maintains the status of the edge pq.
 - The status is *branch* if the edge is known to be in the MST, *reject* if the edge is known not to be in the MST, and *basic* otherwise.
 - father_p: For each process *p* in the fragment, father_p is the edge leading to the core edge of the fragment.
 - state_p: State of node *p* is *find* if *p* is currently engaged in the fragment's search for the lowest-weight outgoing edge and *found* otherwise. Initially it is in state *sleep*.

```
var state<sub>p</sub> : (sleep, find, found) ;
```

```
statch_{p}[q] : (basic, branch, reject) for each q \in Neigh_{p};

name_{p}, bestwt_{p} : real ;

level_{p} : integer ;

testch_{p}, bestch_{p}, father_{p} : Neigh_{p} ;

rec_{p} : integer;
```

(1) As the first action of each process, the algorithm must be initialized:

begin let pq be the channel of p with smallest weight ;

```
statch_p[q] := branch ; level_p := 0 ;
state_p := found ; rec_p := 0 ;
```

```
send (connect, 0) to q
```

```
end
```

INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

(2) Upon receipt of $\langle \text{connect}, L \rangle$ from q:

```
begin if L < level<sup>p</sup> then (* Combine with rule A *)
```

```
begin statch<sub>p</sub>[q] := branch ;
send (initiate, level<sub>p</sub>, name<sub>p</sub>, state<sub>p</sub>) to q
```

```
end
```

```
else if statch_p[q] = basic
then (* Rule C *) process the message later
else (* Rule B *) send (initiate, level_p + 1, \omega(pq), find) to q
```

end

```
(3) Upon receipt of \langle \text{initiate, L, F, S} \rangle from q:
```

```
begin level_p := L; name_p := F; state_p := S; father_p := q;

bestch_p := udef; bestwt_p := \infty;

forall r \in Neigh_p: statch_p[r] = branch \land r \neq q do

send \langle initiate, L, F, S \rangle to r;

if state_p = find then begin rec_p := 0; test end
```

end

Testing the edges

- To find its lowest-weight outgoing edge, node p inspects its outgoing edges in increasing order of weight.
- To inspect edge pq, p sends a (test, level_p, name_p) message to q and waits for an answer
 - A $\langle reject \rangle$ message is sent by process *q* if *q* finds that *p*'s fragment name is the same as *q*'s fragment name. On receiving the $\langle reject \rangle$ message, *p* continues its local search.
 - If the fragment name differs q sends an (accept) message.
 - The processing of a $\langle test, L, F \rangle$ message is deferred by q if L > $level_q$ because pand q may actually belong to the same fragment, but the $\langle initiate, L, F, S \rangle$ message has not yet reached q

A simple optimization

- To inspect edge pq, p sends a (test, level_p, name_p) message to q and waits for an answer
 - A (reject) message is sent by process q if q finds that p's fragment name is the same as q's fragment name.
 - If the edge pq was just used by q to send a (test, L, F) message then p
 will know (in a symmetrical way) that the edge pq is to be rejected. In
 this case, the (reject) message need not be sent by q.

(4) procedure test

```
begin if \exists q \in Neigh_p: statch<sub>p</sub>[q] = basic then
        begin testch_p := q with
                      statch<sub>p</sub>[q] = basic and \omega(pq) minimal;
                 send \langle \text{test}, \text{level}_p, \text{name}_p \rangle to \text{testch}_p
         end
       else begin testch<sub>p</sub> := udef ; report end
end
```

```
(5) Upon receipt of \langle \text{test}, L, F \rangle from q:
```

```
begin if L > level<sub>p</sub> then (* Answer must wait *)
```

```
process the message later
```

```
else if F = name<sub>p</sub> then (* internal edge *)
```

```
begin if statch_p[q] = basic then statch_p[q] := reject;
```

```
if q \neq testch_p
```

```
then send \langle reject \rangle to q
```

else test

end

else send $\langle accept \rangle$ to q

end INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

```
Upon receipt of (accept) from q:
(6)
     begin testch_p := udef;
               if \omega(pq) < best_p
                 then begin bestwt_p := \omega(pq);
                          bestch_p := q
                        end;
                report
     end
(7) Upon receipt of \langle reject \rangle from q:
     begin if statch_p[q] = basic then statch_p[q] := reject;
        test
```

end

```
INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR
```

Reporting the lowest-weight outgoing edge

- The lowest-weight outgoing edge is reported for each sub tree using (report, *w*) messages
 - Node *p* counts the number of (report, ω) messages it receives, using the variable rec_p .
- The (report, *w*) messages are sent in the direction of the core edge by each process.
 - The messages of the two core nodes cross on the edge; both receive the message from their father
 - When this happens, the algorithm terminates if no outgoing edge was reported. Otherwise a $\langle connect, L \rangle$ message must be sent through this edge.

Reorientation of the tree

- If an outgoing edge was reported, the lowest-weight outgoing edge is found by following the *bestch* pointer in each node, starting from the core node on whose side the best edge was reported
- A (connect, L) message must be sent through this edge, and all father pointers in the fragment must point in this direction
 - This is done by sending a (changeroot) message.
 - When the (changeroot) message arrives at the node incident to the lowest-weight outgoing edge, this node sends a (connect, L) message via the lowest-weight outgoing edge.

(8) procedure *report*:

```
begin if rec_p = \#\{q : statch_p[q] = branch \land q \neq father_p\}
and testch_p = udef then
begin state_p := found;
send \langle report, bestwt_p \rangle to father_p
end
```

```
(9) Upon receipt of (report, \omega) from q:
     begin if q \neq father_p
        then (* reply for initiate message *)
          begin if \omega < bestwt_p then
                begin bestwt_p := \omega; bestch_p := q end;
                rec_p := rec_p + 1; report
          end
        else (* pq is the core edge *)
          if state<sub>p</sub> = find
                then process this message later
                else if \omega > bestwt_p
                        then changeroot
                        else if \omega = bestwt<sub>p</sub> = \infty then stop
      end
```

INDIAN INSTITUTE OF TECHNOLOGY KHARAGPUR

(10) procedure *changeroot*:

```
begin if statch_p[bestch_p] = branch
then send (changeroot) to bestch_p
else begin send (connect, level_p) to bestch_p;
statch_p[bestch_p] := branch
end
end
```

(11) Upon receipt of ⟨changeroot⟩:begin *changeroot* end

Complexity

The Gallager-Humblet-Spira algorithm computes the minimal spanning tree using at most 5N logN + 2|E| messages

- Each edge is rejected at most once and this requires two messages (test and reject). This accounts for at most 2|E| messages.
- At any level, a node receives at most one initiate and one accept message, and sends at most one report, one changeroot *or* connect message, and one test message not leading to a rejection. For logN levels, this accounts for a total of 5N logN messages.