

BILKENT UNIVERSITY

CS 442

DISTRIBUTED SYSTEMS

Cooperative Routing Protocol

Halil ATMACA
21302459

Burak GÖK
21302025

Enes VAROL
21604086

Instructor:
Buğra GEDİK
Assoc. Prof. of Computer Eng.

Submitted on May 27, 2019
Revised on Sep 25, 2019



1 Introduction

Communication networks are generally not fully connected, which necessitates routing of messages sent among nodes. Routing requires each node in a network to select a next-hop node for every destination. While the selection is trivial for neighbors, it involves the computation of shortest distance for non-neighboring nodes.

If the topology is known prior to the computation, which is the case in link-state routing protocols, a single-source shortest path algorithm, such as the Dijkstra's algorithm, can be used to form the routing table. In such protocols, each node sends its link-state advertisement, which is the collection of its neighbors along with their distances, to every node in the network via flooding. Once a node has the complete set of link-state advertisements, it forms the network graph and its routing table. Since this is a cumulative process, the obtained results may not completely reflect the current state of the network; or, they may contain contradictory information. Usually sequence-numbered messages are used to overcome outdated or contradictory information. Link-state routing protocols are very restrictive for highly dynamic networks, due to their computational complexity and the number of messages they generate in result of flooding.

There is another family of routing protocols, which does not require obtaining the complete network image in order to compute shortest distances. Instead, in such protocols, each node stores the distance to every other node along with the next-hop node associated with it. The array of distances to other nodes is known as the distance vector. When a distance vector is manipulated, it is sent to all neighbors. When a node receives a distance vector from a neighbor, it updates its own vector to reflect changes in network topology and link weights. Since the distance vector of a node functions like its routing table along with the next-hop nodes, the comparison of distance vectors is the only computation done by nodes. For such comparisons, a distributed variant of the Bellman-Ford algorithm is used. The algorithm is slower than the Dijkstra's algorithm, but the former can be decentralized unlike the latter. Distance-vector routing protocols are simpler and has less message overhead compared to link-state routing protocols. However, since they do not try to form the complete network graph, they cannot test the consistency of their distance vectors, which can lead to routing loops and the widely known count-to-infinity problem.

In this project, we try to devise a distance-vector routing protocol that avoids routing loops and the count-to-infinity problem and also has the optimal number of message exchanges. Our protocol is purely theoretical and requires extra work to be implemented in real communication networks. The protocol has the following limitations or relaxations.

1. Network links must have non-negative weights.
2. Reliable message exchange protocol is required such as TCP.
3. Neighboring nodes may have asymmetric weights for the same link.
4. The triangle inequality does not have to be satisfied among network links.
5. Any distance metric can be used as long as from two paths, the more preferable one has the smaller distance.

2 Protocol

The protocol is based on the Bellman-Ford algorithm as usual. It differs from similar protocols by introducing two properties for nodes: *smartness* and *helpfulness*. A node is smart if it always preserves the consistency of its internal state and does not lose any information that might be helpful while forming the routing table. A node is helpful if it does not spread any information that may mislead its neighbors and informs them if their routing tables are suboptimal.

The protocol utilizes a reachability protocol, to obtain information about neighboring nodes such as link weight changes and mobility events. This is done by sending echo requests and collecting echo replies at each node. The time measured between a request-reply pair is used to assign a weight to the respective link. If no reply is received after the request within the timeout period, the link weight is set to infinity. In that case, the link is said to be congested and for a while, echo requests are continued to send. If the respective neighbor does not respond for an unusual amount of time, it is forgotten. Receiving an echo request from an unknown or non-neighboring node means that a new link is discovered.

When a node joins the network for the first time or a disconnected node becomes connected again, it notifies its neighbors by sending echo requests during a usual echo process. When a node A receives an echo request from an unknown node B, it does not add node B to the distance vector because the distance to node B is not measured yet. After node A initiates an echo process, the distance to node B can be measured after receiving the echo reply. However, node B may send a distance vector to node A before sending an echo reply since it has node A in its distance vector. It is impossible to compare a distance vector without knowing the distance to its sender. That's why the protocol enforces node B to put its distance to node A in the distance vector it sends to node A. This is called *the reported distance* and is only used in such a situation. For all cases, the distance obtained by echoing is preferred.

Since the triangle inequality does not have to be satisfied among links, a neighboring node may have a next-hop node other than itself—or NIL, which is denoted as \emptyset . This causes such a node to have a distance that is different from its link weight. Thus, the node distance and the link weight should be stored separately to satisfy the second requirement of smartness property because there may be an update, in the future, which causes switching between two quantities. We will use the following symbols in the rest of the report for the sake of simplicity and formalism.

L	Links (All neighbors)
$W : L \rightarrow \mathbb{R}^+$	Link weight
$N : \{x \in L \mid W(x) \neq \infty\}$	Responsive neighbors
$T : L \rightarrow \mathbb{R}^+$	Time of the last echo request/reply received
$T(\text{REQ})$	Time of the last echo request sent
A	All nodes in the network
K	Known nodes
$P : K \rightarrow N \cup \{\emptyset\}$	Predecessor (Next-hop node)
$D : K \rightarrow \mathbb{R}^+$	Node distance
$U : R \rightarrow \mathbb{R}^+, R \subseteq A$	Distance vector received
M	Myself
S	Sender

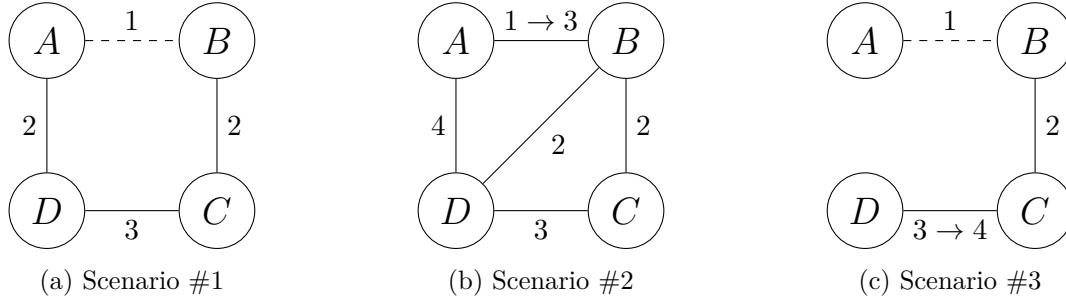


Fig. 1. Topology change scenarios

2.1 Smartness

As the first requirement of smartness, nodes should preserve the consistency of their internal state, i.e. distance vector, after an update. For example, in Fig. 1a, assume that the network is stabilized. After the link between node A and node B is broken and node A realizes that node B is unreachable, it sets $D(B)$ to infinity. However, node C is reached via node B , so $D(C)$ should also be set to infinity. In general, when the distance to a neighbor is updated, the distances to all descendants of that neighbor should also be updated. Such an update is denoted as $D(x) \leftarrow y$.

definition $D(x) \leftarrow y$:

- 1 **for each** $k \in K$ **s.t.** $P(k) = x$ **do**
- 2 $D(k) += y - D(x)$
- 3 $D(x) = y$

When a neighbor is decided to reach indirectly as a result of an update, which can only occur while processing a distance vector received, the predecessors of all descendants of that neighbor should also be updated. In order not to increase the complexity of link weight updates, such an update is denoted differently as $D(x) \leftarrow y$.

definition $D(x) \leftarrow y$:

- 1 **for each** $k \in K$ **s.t.** $P(k) = x$ **do**
- 2 $P(k) = P(x) \neq \emptyset ? P(x) : x$
- 3 $D(k) += y - D(x)$
- 4 $D(x) = y$

As stated before, for the second requirement of smartness, nodes should keep information that might be useful while forming the routing table. Assume that the network depicted in Fig. 1b is stabilized. Node A reaches node D via node B . After the weight of the link between node A and node B becomes 3, node A should switch to the link between node D and itself to reach node D . In general, if the link weight of an indirectly reached neighbor is smaller than its distance, the link should be preferred. If the link weight and node distances are equal, the direct link is still preferred for the sake of stability.

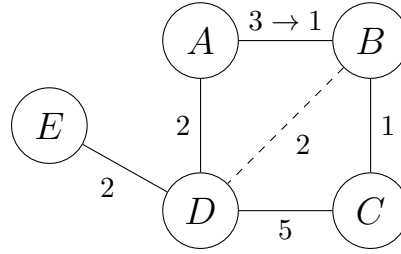


Fig. 2. A topology change scenario that results in a suboptimal state

Pitfalls

When a neighbor is decided to reach indirectly as a result of an update, its descendants might be processed twice: once due to the first requirement of smartness and once due to their probable inclusion in the distance vector received. If the sender node has an optimal state, the distances that are processed for the same node will be equal and the order of processing does not matter. Otherwise, the smaller one should be selected and in some cases, the update may result in a suboptimal state regardless of the order of processing. In Fig. 2, assume that the link between node B and node D is not established yet and the network is stabilized. In this topology, node C reaches node E via node D . If the aforementioned link is established and the weight of the link between node A and node B becomes 1 simultaneously, node C will receive an update from node B , which claims that $D(D) = 2$ and $D(E) = 5$. Two possible scenarios for the order of processing are as follows.

1. If the former is processed first, node D will be decided to reach via node B and consequently the distance of node E will be set to 5. When the latter is processed, the distance of node E will be updated again and set to 6 since it is now reached via node B .
2. If the latter is processed first, node E will be decided to reach via node B and its distance will be set to 6. When the former is processed, the distance of node E will not be updated since it is not reached via node D anymore.

The root cause of this problem is to compute the candidate paths to reach non-neighboring nodes using the modified distance vector instead of the original one. Thus, it is possible to avoid this pitfall by keeping the original distance vector. In this way, there will be two distance vectors: the backup vector, using which the candidate paths are computed, and the resultant vector, with which the candidate paths are compared. After all nodes are processed, the resultant vector will have the shortest paths. This solution requires $D(x) \leftarrow y$ to be modified as the following, where K' , P' and D' respectively hold the initially known nodes, predecessors of and distances to them.

definition $D(x) \leftarrow y$:

```

1  for each  $k \in K'$  s.t.  $P'(k) = x$  do
2       $D'' = D'(k) + (y - D'(x))$ 
3      if  $D'' < D(k)$  then
4           $P(k) = P(x) \neq \emptyset ? P(x) : x$ 
5           $D(k) = D''$ 
6   $D(x) = y$ 
```

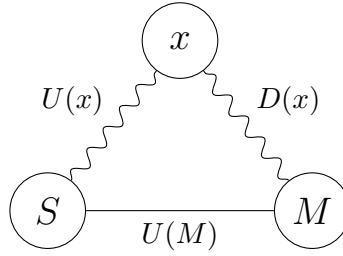


Fig. 3. Visualization of an update

2.2 Helpfulness

As the first requirement of helpfulness, nodes should not spread information that can be misleading. In Fig. 1c, if the link between node A and node B is broken and the weight of the link between node C and node D becomes 4 simultaneously, node B will receive an update from node C , which claims that $D(A) = 3$, after setting node A as unreachable. Since node B has no way of knowing that node C actually reaches node A via itself, it will apply the update. This will eventually lead to the count-to-infinity problem. If node A did not become disconnected after the link is broken; for example, if there was a link between node A and node D as in Fig. 1a, the network would still converge. However, the convergence time and the message overhead would increase. To overcome these problems, when a node sends its distance vector to another node, only the distances of the nodes that are not reached via the destination node are sent. This is expressed more formally as "send $\{(k, D(k)) \mid k \in K, P(k) \neq n\}$ to each $n \in N$."

As the second requirement of helpfulness, nodes should inform their neighbors if their routing tables are suboptimal. Since one of the main objectives of the protocol is to decrease the number of messages exchanged, distance vectors are only broadcast if there is an update. Furthermore, when they are broadcast, the message-originator neighbor that causes the update is skipped because it cannot learn new information from an update it originated. Because of these two reasons, sometimes nodes may stuck in a suboptimal, or possibly faulty state.

In Fig. 1a, assume that node B realizes that the link between node A and itself is broken before node A . Node B will set $D(A)$ to infinity and send its distance vector to node C . Since node C reaches node A via node B , it will also set $D(A)$ to infinity and send its distance vector to node D . However, node D will not update its distance vector since it reaches node A via a direct link. After node A realizes that node B is unreachable, it will set $D(B)$ to infinity and sends its distance vector to node D . Since node D reaches node B via node A , it will also set $D(B)$ to infinity and send its distance vector to node C . Since node C reaches node B via a direct link, it will not update $D(B)$. However, it will update $D(A)$ since it realizes that node A is reachable via node D . Note that node A and node D still believe that node B is unreachable, which cannot be corrected since the network is now stabilized.

A node's state can be corrected if its neighbors detect possible improvements in its distance vector and informs it. Fig. 3 visualizes an update, where S and M are the sender and receiver nodes respectively. If $U(M) + D(x) < U(x)$, node S can reach node x via node M , which is shorter than the path it currently uses.

The complete protocol is as follows. The following operators are defined to manipulate sets.

$$\begin{aligned} x += y &\equiv x = x \cup \mathbb{S}(y) \\ x -= y &\equiv x = x \setminus \mathbb{S}(y) \end{aligned} \quad \text{where} \quad \mathbb{S}(x) = \begin{cases} x & \text{if } x \text{ is a set,} \\ \{x\} & \text{otherwise.} \end{cases}$$

procedure RECEIVE(Echo Request)

```

1  if  $S \notin L$  then
2       $L += S$ 
3       $W(S) = \infty$ 
4       $T(S) = \text{NOW}$ 
5      send echo request to  $S$ 
```

procedure RECEIVE(Echo Reply)

```

1  if  $S \notin L$  then
2       $L += S$ 
3       $W(S) = \infty$ 
4       $T(S) = \text{NOW}$ 
```

procedure ECHO()

```

1   $T(\text{REQ}) = \text{NOW}$ 
2  send echo request to each  $x \in L$ 
3  wait for echo timeout

4   $K', P', D' = K, P, D$ 
5  for each  $x \in L$  do
6       $W(x) = T(x) > T(\text{REQ}) ? T(x) - T(\text{REQ}) : \infty$ 
7      if  $x \notin K$  then
8          if  $W(x) \neq \infty$  then
9               $K += x$ 
10              $P(x) = \emptyset$ 
11              $D(x) = W(x)$ 
12         else if  $W(x) \leq D(x)$  or  $P(x) = \emptyset$  then
13              $P(x) = \emptyset$ 
14              $D(x) \leftarrow W(x)$ 
15         if  $T(\text{REQ}) - T(x) \geq \text{link life}$  then
16              $L -= x$ 

17  updated  $\equiv (\exists x \in K) [x \notin K' \text{ or } P'(x) \neq P(x) \text{ or } D'(x) \neq D(x)]$ 
18  if updated then
19      send  $\{(k, D(k)) \mid k \in K, P(k) \neq n\}$  to each  $n \in N$ 
20       $K -= \{k \in K \mid D(k) = \infty\}$ 
```

procedure RECEIVE(U)

```

1   $K', P', D' = K, P, D$ 
2  if  $S \notin L$  or  $W(S) = \infty$  then
3       $L += S$ 
4       $W(S) = U(M)$ 
5  if  $S \notin K$  then
6       $K += S$ 
7       $P(S) = \emptyset$ 
8       $D(S) = U(M)$ 

9  for each  $x \in R \setminus \{M\}$  do
10      $D'' = D(S) + U(x)$ 
11     if  $x \notin K$  then
12         if  $D'' \neq \infty$  then
13              $K += x$ 
14              $P(x) = S$ 
15              $D(x) = D''$ 
16         else if  $D'' < D(x)$  or  $P'(x) = P(x) = S$  then
17              $P(x) = S$ 
18              $D(x) \leftarrow D''$ 
19     for each  $x \in N$  s.t.  $W(x) \leq D(x)$  do
20          $P(x) = \emptyset$ 
21          $D(x) = W(x)$ 

22     updated  $\equiv (\exists x \in K) [x \notin K' \text{ or } P'(x) \neq P(x) \text{ or } D'(x) \neq D(x)]$ 
23     inform  $\equiv (\exists x \in R \setminus \{M\}) [U(M) + D(x) < U(x)]$ 
24     if updated or inform then
25          $N' = \{n \in N \mid (\text{updated and } n \neq S) \text{ or } (\text{inform and } n = S)\}$ 
26         send  $\{(k, D(k)) \mid k \in K, P(k) \neq n\}$  to each  $n \in N'$ 
27          $K -= \{k \in K \mid D(k) = \infty\}$ 

```