

## A APPENDIX A

This appendix provides pseudo-code to some of the sparsification algorithms as a supplement to the paper.

### A.0.1 Rank Degree Sparsifier.

---

#### Algorithm 1 Rank Degree sparsifier

---

```

1: procedure RANKDEGREESPARSIFIER( $G$ )
2:   Input:  $G$ : Graph to sparsify
3:   Input:  $\rho: 0 < \rho \leq 1$  selects top  $\rho * \#$ neighbors for each vertex
4:   Output:  $H$ : Sparsified graph
5:
6:   seeds = [ $u_1, u_2, \dots, u_s$ ]
7:    $\mathcal{V}_H = \emptyset, \mathcal{E}_H = \emptyset$ 
8:   while  $|\mathcal{V}_H| < |\mathcal{V}_G|$  do
9:     new_seeds= $\emptyset$ 
10:    for all  $u \in$  seeds do
11:      neihgs=getNeighborsOf( $u$ )
12:      ranks = {}
13:      for all  $v \in$  neihgs do
14:        ranks[ $v$ ] = getDegreeOf( $v$ )
15:      sort ranks by value
16:      select top  $k=\rho \times \text{len(neighbors)}$ ,  $v_1, \dots, v_k$ 
17:      new_seeds = new_seeds  $\cup [v_1, \dots, v_k]$ 
18:       $\mathcal{E}_H = \mathcal{E}_H \cup [(u, v_1), \dots, (u, v_k)]$ 
19:      seeds=new_seeds
20:    $H = \{\text{Vertex}(\mathcal{V}_G), \text{Edge}(\mathcal{E}_H)\}$ 
21:   return  $H$ 

```

---

### A.0.2 Local Degree Sparsifier.

---

#### Algorithm 2 Local degree score

---

```

1: procedure GETEDGESCORE( $G$ )
2:   Input:  $G$ : Graph to calculate edge scores
3:   Output: Scores: An array of edge scores for each edge
4:
5:   scores = [0, ..., 0]
6:   for all  $v_i \in \mathcal{V}$  do
7:      $d_i = \text{degree}(v_i)$ 
8:     neighbor_degree = {}
9:     for all  $v_j \in \text{Neighbor}(v_i)$  do
10:        $d_j = \text{degree}(v_j)$ 
11:       eid =  $e_{ij}.\text{edgeID}$ 
12:       neighbor_degree[eid] =  $d_j$ 
13:
14:   sort(neighbor_degree)
15:
16:   last_rank, last_degree, num_same = 0
17:   neighbor_rank = {}
18:   for all item  $\in$  neighbor_degree do
19:     eid = item.key
20:      $d_j = \text{item.value}$ 
21:     if  $d_j == \text{last\_degree}$  then
22:       num_same++
23:     else
24:       last_rank += num_same
25:       num_same = 1
26:       last_degree =  $d_j$ 
27:     neighbor_rank[eid] = last_rank
28:
29:   for all item  $\in$  neighbor_degree do
30:     eid = item.key
31:     rank = item.value
32:      $s = 1.0 - \log(\text{rank})/\log(d_j)$ 
33:     // score for an edge can be updated multiple times, take the max score
34:     scores[eid] = max(scores[eid], s)
35:   return scores

```

---

### A.0.3 *t*-Spanner.

---

#### Algorithm 3 Greedy algorithm for t-spanner construction

---

```

1: procedure CONSTRUCT T-SPANNER( $G, t$ )
2:   Input:  $G$ : Original graph
3:   Input:  $t$ : stretch factor, must be an odd number  $> 1$ 
4:   Output:  $H$ : t-spanner graph
5:
6:    $H \leftarrow (V, \emptyset)$                                          ▷ Init  $H$  to have the same set of vertices, and no edges
7:   for all  $e_{uv} \in E$  in non-decreasing order do
8:     if  $d_H(u, v) > tw(u, v)$  then
9:       add  $e_{uv}$  to  $H$ 
10:  return  $H$ 

```

---

### A.0.4 Forest Fire.

The *Forest Fire* model can be described more formally as follows (modified from [?]):

- (1) A new vertex  $u$  chooses an existing vertex  $v$  uniformly at random, and forms an edge to it.
- (2) Two random numbers  $x$  and  $y$  are generated geometrically distributed with means  $p/(1-p)$  and  $rp/(1-rp)$ , where  $p$  is the forward burning probability, and  $r$  is the backward burning ratio.
- (3) Vertex  $v$  selects  $x$  outgoing edges and  $y$  incoming edges that are not visited yet, if there are not enough unvisited edges, select all. For undirected edges, every edge can be both an outgoing and incoming edge. Let  $w_1, w_2, \dots, w_{x+y}$  denote the other end of the selected edges.
- (4) Vertex  $u$  forms edges to the  $w_1, w_2, \dots, w_{x+y}$ .
- (5) Repeat (3) and (4) recursively to each of the  $w_1, w_2, \dots, w_{x+y}$ , until no edge can be added.

---

#### Algorithm 4 Forest Fire Score

---

```

1: procedure GETEDGESCORE( $G$ )
2:   Input:  $G$ : Graph to calculate edge scores
3:   Input:  $bp$ : The probability a neighbor vertex is burnt, from 0.0 to 1.0
4:   Input: targetBurnRatio: In total targetBurnRatio * m edges will be burnt
5:   Output: Scores: An array of edge scores for each edge
6:
7:   burnt_count = 0                                         ▷ keep track of total number of burnt edges
8:   scores = [0, ..., 0]                                     ▷ scores is an array of length #edges
9:   burnt = [0, ..., 0]                                     ▷ burnt is an array of length #edges
10:
11:  while burnt_count < targetBurnRatio * number_of_edges( $G$ ) do
12:    visited = [false, ..., false]                           ▷ visited is an array of length #vertices
13:    vertexQ = []                                         ▷ a queue for vertex to be visited
14:    vertexQ.add(randomVertex( $G$ ))                         ▷ pick a random starting vertex
15:
16:    while vertexQ is not empty do
17:      u = vertexQ.pop()                                    ▷ r is a random float from 0.0 to 1.0
18:      visited[u.id] = True
19:      neigns = getAllUnvisitedVertices()
20:      while neigns is not empty do
21:        r = randNum()
22:        if  $r \leq bp$  then
23:          break
24:        v = pickRandom(neigns)
25:        remove(neigns, v)
26:        vertexQ.add(v)
27:        eid = getEdgeID(u, v)
28:        burnt_count++                                     ▷ decides to burn the vertex, not propagate further
29:        max_burnt = max(burnt)                            ▷ pick a random vertex to propagate fire
30:        scores[eid] = burnt / max_burnt                  ▷ pick a random vertex to propagate fire
31:   return scores                                         ▷ normalize burnt_count to be the scores

```

---

### A.0.5 Similarity-based sparsifiers.

---

### Algorithm 5 G-Spar

---

```

1: procedure G-SPAR(G)
2:   Input: G: Graph to calculate edge scores
3:   Output: H: G-Spar sparsified graph
4:
5:   scores = {}
6:   for all e ∈ E do
7:     eid = e.id
8:     score = JaccardScore(e)
9:     scores[eid] = score
10:    sort scores by value
11:    pick top s% edges to form H
12:   return H

```

---



---

### Algorithm 6 L-Spar

---

```

1: procedure L-SPAR(G)
2:   Input: G: Graph to calculate edge scores
3:   Input: c: Exponent parameter
4:   Output: H: L-Spar sparsified graph
5:
6:   for all v ∈ V do
7:     d = degreeOf(v)
8:     E'=getEdgesOf(v)
9:     scores = {}
10:    for all e ∈ E' do
11:      eid = e.id
12:      score = JaccardScore(e)
13:      scores[eid] = score
14:    sort scores by value
15:    add top dc edges to H
16:   return H

```

---



---

### Algorithm 7 Edge Triangle Count

---

```

1: procedure EDGETRIANGLECOUNT(G)
2:   Input: G: Graph to calculate triangle edge scores
3:   Output: triangle_count: Triangle count for each edge
4:
5:   triangle_count = [0, ..., 0]                                ▷ array of length #edges
6:   incident_triangle_count = [None, ..., None]                ▷ array of length #vertices
7:
8:   for all u ∈ V do                                         ▷ first vertex in triangle
9:     for all v ∈ getNeighborsOf(u) do                         ▷ mark all neighboring vertices not None
10:      incident_triangle_count[v] = 0
11:      for all v ∈ getNeighborsOf(u) do                         ▷ second vertex in triangle
12:        for all w ∈ getNeighborsOf(v) do                         ▷ third vertex in triangle
13:          if incident_triangle_count[w] is not None then           ▷ triangle found
14:            // count triangles to the vertices first, each triangle is counted 3 times
15:            if u ≥ v then
16:              incident_triangle_count[v]++
17:            if u ≥ w then
18:              incident_triangle_count[w]++
19:
20:  // add local triangle count to global, reset local triangle count
21:  for all v ∈ getNeighborsOf(u) do
22:    eid = getEdgeId(u, v)
23:    if incident_triangle_count[v] > 0 then
24:      triangle_count[eid] += incident_triangle_count[v]
25:      incident_triangle_count[v] = None
26:   return triangle_count

```

---

The Edge Triangle Count is not a standalone sparsification algorithm. We list it here because it is used in the calculation of the local similarity score and the SCAN structural similarity score.

---

**Algorithm 8** Local similarity score

---

```

1: procedure LOCALSIMILARITYSCORE( $G$ )
2:   Input:  $G$ : Graph to calculate triangle edge scores
3:   Output: scores: Local similarity scores for each edge
4:
5:   scores = [0, ..., 0]                                     ▷ array of length #edges
6:
7:   triangle_count = EdgeTriangleCount( $G$ )
8:   for all  $u \in V$  do
9:     neighbors_sims = {}
10:     $d_u$  = getDegreeOf( $u$ )
11:    for all  $v \in \text{getNeighborsOf}(u)$  do
12:       $d_v$  = getDegreeOf( $v$ )
13:      eid = getEdgeId( $u, v$ )
14:      sim = triangle_count[eid]/( $d_u + d_v - \text{triangle\_count}[eid]$ )
15:      scores[eid] = max(scores[eid], sim)
16:   return scores

```

---



---

**Algorithm 9** SCAN Structural Similarity Score

---

```

1: procedure SCANSTRUCTURALSIMILARITYSCORE( $G$ )
2:   Input:  $G$ : Graph to calculate triangle edge scores
3:   Output: scores: SCAN structural similarity scores for each edge
4:
5:   scores = [0, ..., 0]                                     ▷ array of length #edges
6:
7:   triangle_count = EdgeTriangleCount( $G$ )
8:   for all  $u \in V$  do
9:     neighbors_sims = {}
10:     $d_u$  = getDegreeOf( $u$ )
11:    for all  $v \in \text{getNeighborsOf}(u)$  do
12:       $d_v$  = getDegreeOf( $v$ )
13:      eid = getEdgeId( $u, v$ )
14:      sim = ( $\text{triangle\_count}[eid] + 1$ ) /  $\sqrt{(d_u + 1) * (d_v + 1)}$ 
15:      scores[eid] = sim
16:   return scores

```

---

#### A.0.6 Effective Resistance (ER) Sparsifier.

We briefly summarize the derivation of the effective resistance. Interested readers should refer to [?] for more details.

We first define the following notations:

$\mathbb{R}$ : real number.

$G$ : Input Graph, in this write-up,  $G$  must be symmetrical (undirected).

$|\mathcal{V}|$ : Number of Vertices in  $G$ .

$|\mathcal{E}|$ : Number of Edges in  $G$ .

$A \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ , Adjacency Matrix of  $G$ .

$D \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ , Degree Matrix of  $G$ , where  $i^{th}$  diagonal entry is the degree of  $i^{th}$  vertex, if the graph is weighted, then it's the sum of all edge weights related to vertex  $i$ .

$L \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ , Laplacian Matrix of  $G$ ,  $L = D - A$ .

$B$ : Incidence Matrix,  $\in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{V}|}$ . Each row in  $B$  represents an edge, where the head vertex is -1, the tail vertex is 1, and all others are 0s. The head and tail of an undirected edge are randomly assigned.

$W$ : Weight Matrix,  $\in \mathbb{R}^{|\mathcal{E}| \times |\mathcal{E}|}$ , is a diagonal matrix, and each diagonal entry represents an edge weight. If the graph is unweighted, then  $W$  becomes an Identity Matrix  $I$ .

$\chi_u$ : A unit vector of length  $|\mathcal{V}|$ , where only the  $u^{th}$  element is 1, others are 0s.

$R_{uv}$ : The effective resistance of edge  $uv$ .

Now we start the derivation of the effective resistance:

$$L = B^T W B \quad (\text{proof omitted}) \tag{1a}$$

According to Kirchhoff's law, the current flow in is always the same as the current flow out of the vertex,

$$\mathbf{B}^T \mathbf{i} = \mathbf{c}_{ext} \quad (1b)$$

According to Ohm's law,

$$\mathbf{i} = \mathbf{W}\mathbf{B}\mathbf{v} \quad (1c)$$

Combining eq. (1a), (1b), and (1c),

$$\mathbf{B}^T \mathbf{W}\mathbf{B}\mathbf{v} = \mathbf{L}\mathbf{v} = \mathbf{c}_{ext} \quad (1d)$$

Let  $\mathbf{L}^+$  be the pseudo-inverse of  $\mathbf{L}$ , because Laplacian matrix is positive semi-definite, and doesn't have an inverse

$$\mathbf{v} = \mathbf{L}^+ \mathbf{c}_{ext} \quad (1e)$$

Now set  $\mathbf{c}_{ext} = \chi_u - \chi_v$ , then eq. (1e) can be written as

$$\mathbf{v} = \mathbf{L}^+ (\chi_u - \chi_v) \quad (1f)$$

Multiply both sides by  $(\chi_u - \chi_v)^T$ ,

$$(\chi_u - \chi_v)^T \mathbf{v} = (\chi_u - \chi_v)^T \mathbf{L}^+ (\chi_u - \chi_v) \quad (1g)$$

Notice that  $(\chi_u - \chi_v)$  is equivalent to the transpose of  $i^{th}$  row in  $\mathbf{B}$ , denoted by  $\mathbf{B}[i]$ , thus,

$$\mathbf{B}[i]^T \mathbf{v} = \mathbf{B}[i] \mathbf{L}^+ \mathbf{B}[i]^T \quad (1h)$$

Eq. (1h) applies to every  $i$ , thus can be generalized to

$$\mathbf{B}^T \mathbf{v} = \mathbf{B} \mathbf{L}^+ \mathbf{B}^T \quad (1i)$$

The l.h.s. of eq. (1g) is the voltage difference between  $u$  and  $v$ , which can be used to represent the effective resistance of the edge connecting  $u$  and  $v$ . Thus, the effective resistance is defined as

$$\begin{aligned} R_{uv} &= (\chi_u - \chi_v)^T \mathbf{L}^+ (\chi_u - \chi_v) \\ &= (\chi_u - \chi_v)^T \mathbf{L}^+ \mathbf{L} \mathbf{L}^+ (\chi_u - \chi_v) \\ &= (\chi_u - \chi_v)^T \mathbf{L}^+ \mathbf{B}^T \mathbf{W} \mathbf{B} \mathbf{L}^+ (\chi_u - \chi_v) \\ &= ((\chi_u - \chi_v)^T \mathbf{L}^+ \mathbf{B}^T \mathbf{W}^{1/2}) (\mathbf{W}^{1/2} \mathbf{B} \mathbf{L}^+ (\chi_u - \chi_v)) \\ &= \|\mathbf{W}^{1/2} \mathbf{B} \mathbf{L}^+ (\chi_u - \chi_v)\|_2^2 \end{aligned} \quad (1j)$$

## B APPENDIX B

Due to the page limit, we only showed a subset of the results in the paper. This appendix presents the full results generated in this work. We show one dataset on each page, and the sub-captions note which metric each subgraph is measuring. Some figures are missing, and we explain why they are missing in such cases. There are three reasons why certain figures are missing, 1) Some metrics are supported for directed graphs, they are Clustering F1 Similarity, Number of Communities, and Modularity; 2) Some experiments couldn't finish within 24 hours, especially on time-consuming metrics like Eigenvector Centrality and large graphs; 3) Some experiments run out of memory and triggered OOM kill by the OS.

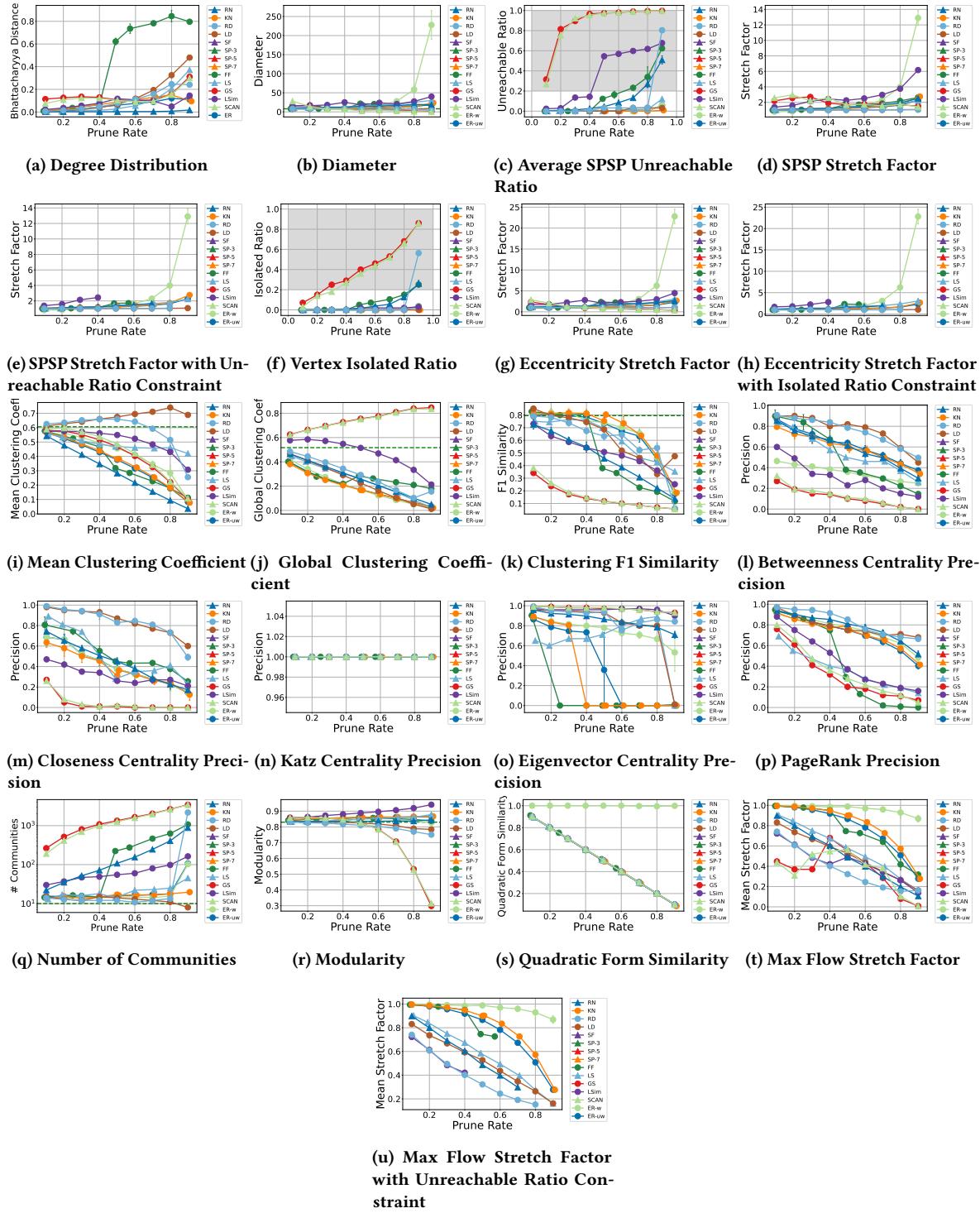


Figure 1: Metric Evaluation on ego-Facebook

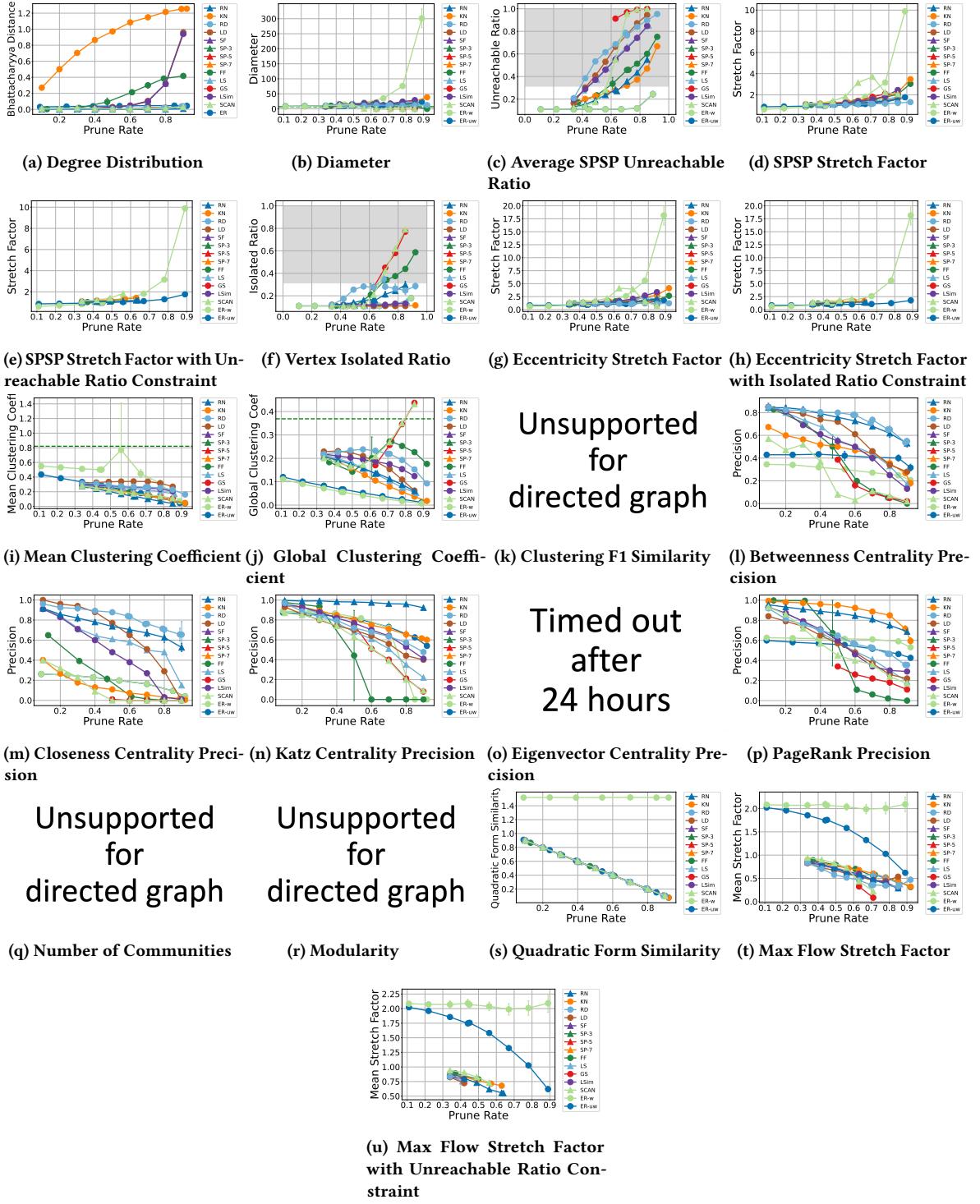


Figure 2: Metric Evaluation on ego-Twitter

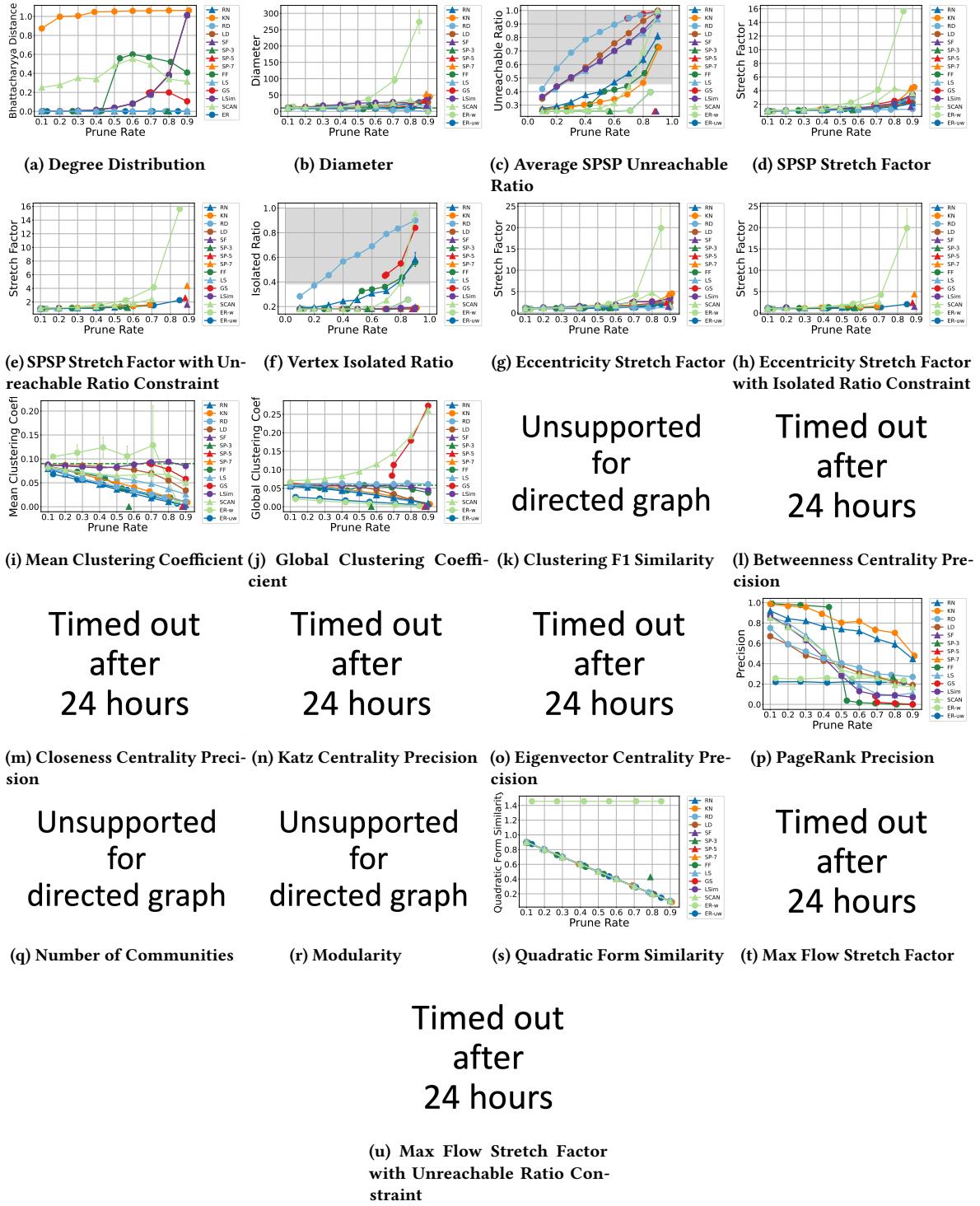


Figure 3: Metric Evaluation on soc-Pokec

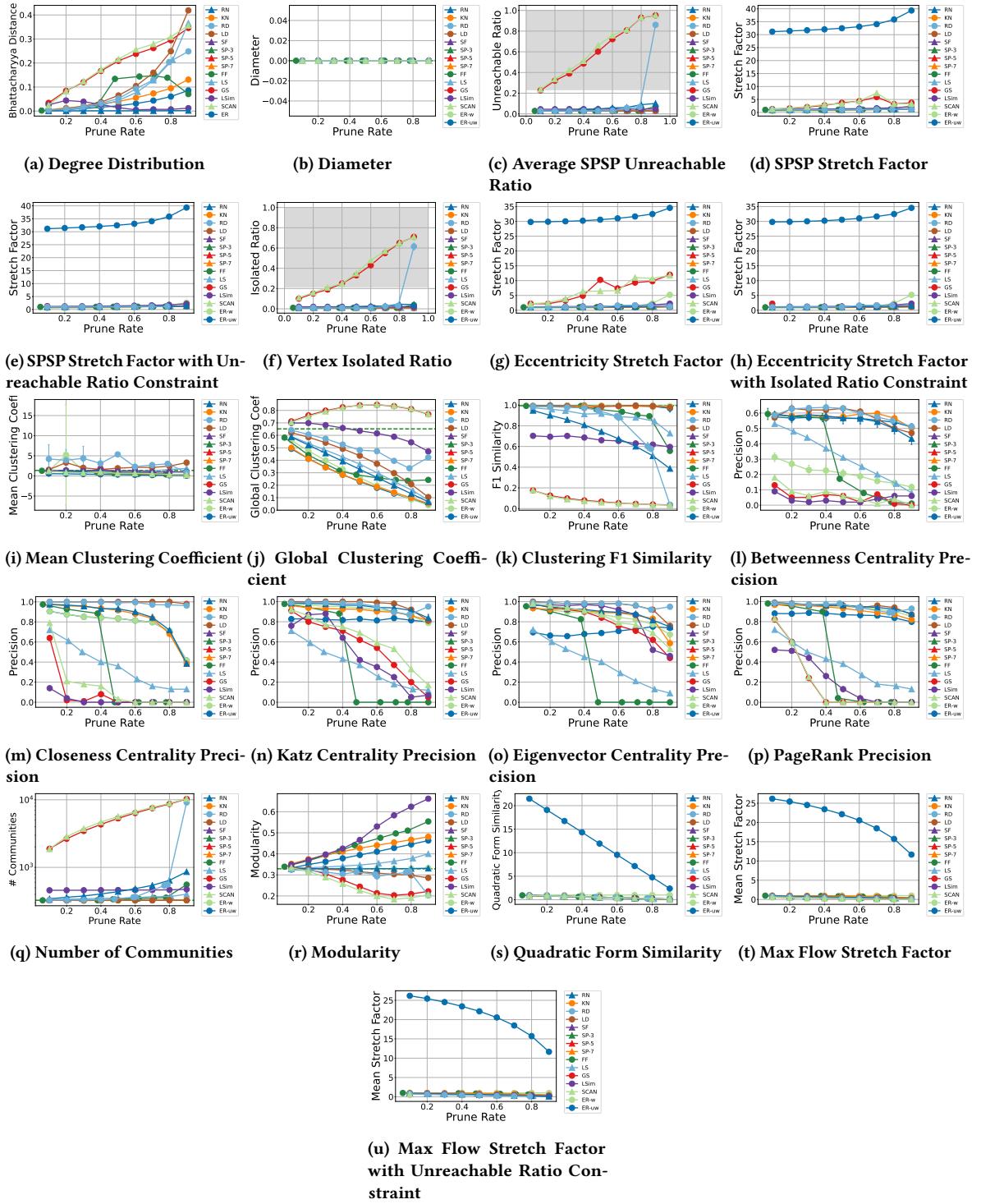


Figure 4: Metric Evaluation on `human_gene2`

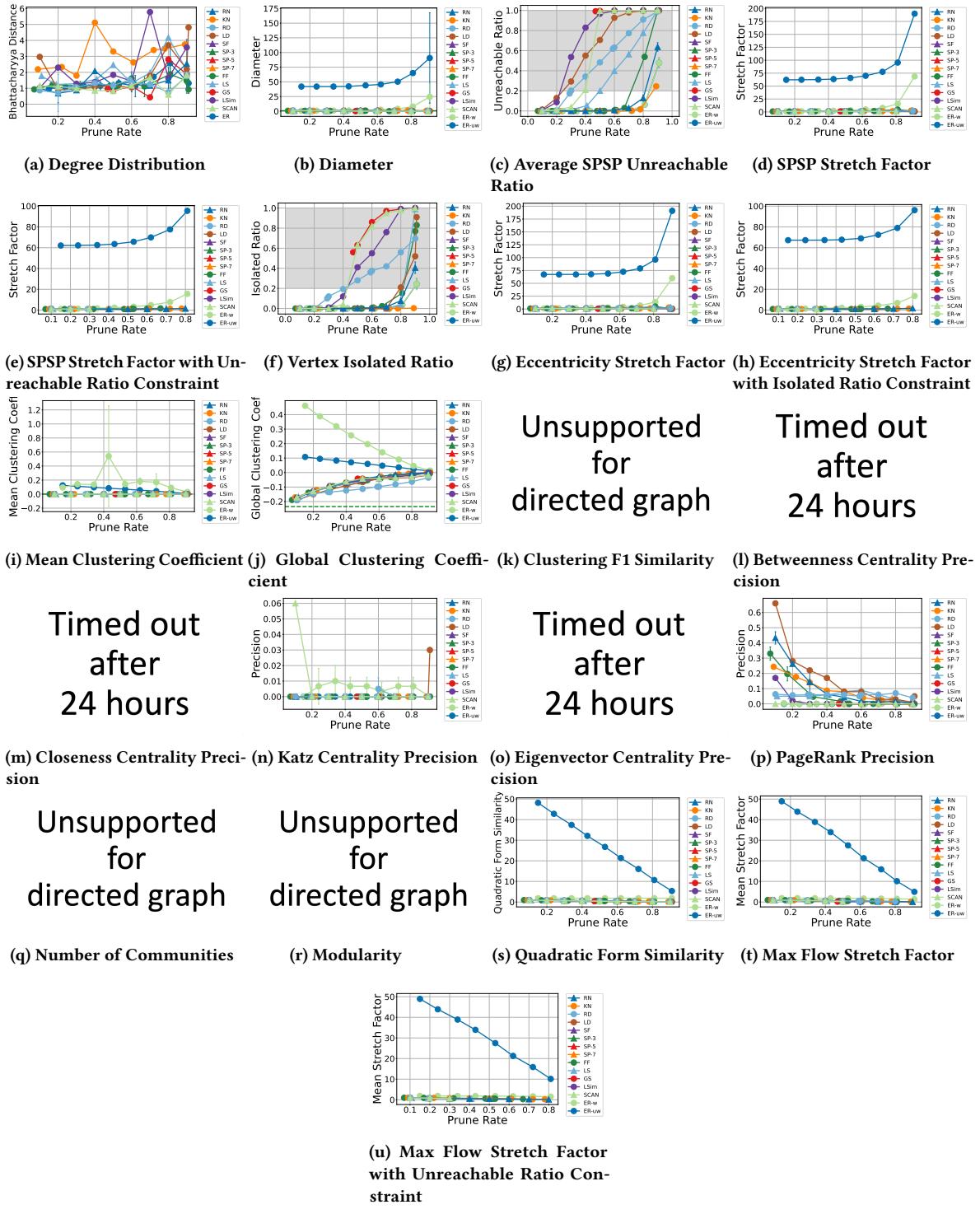


Figure 5: Metric Evaluation on cage14

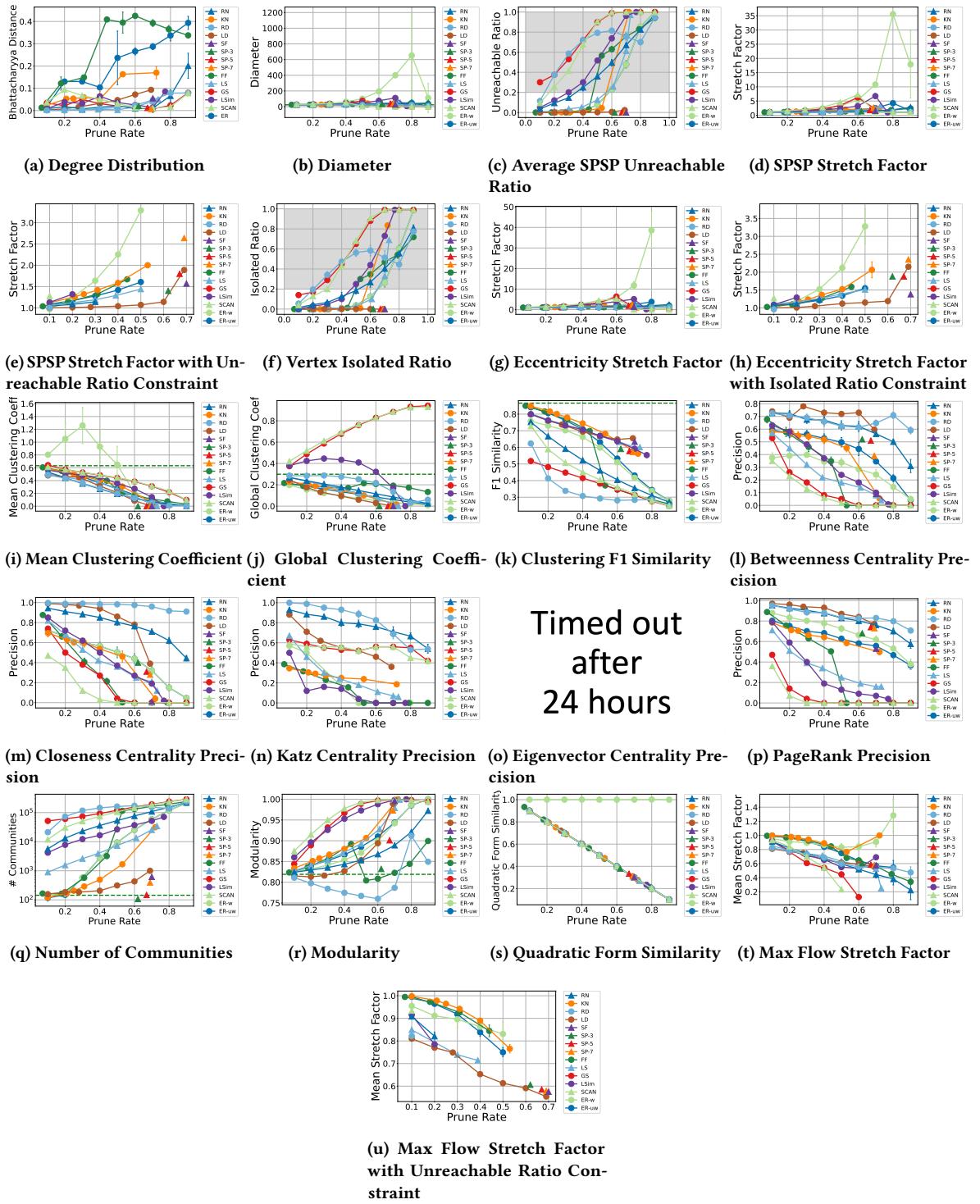


Figure 6: Metric Evaluation on com-DBLP

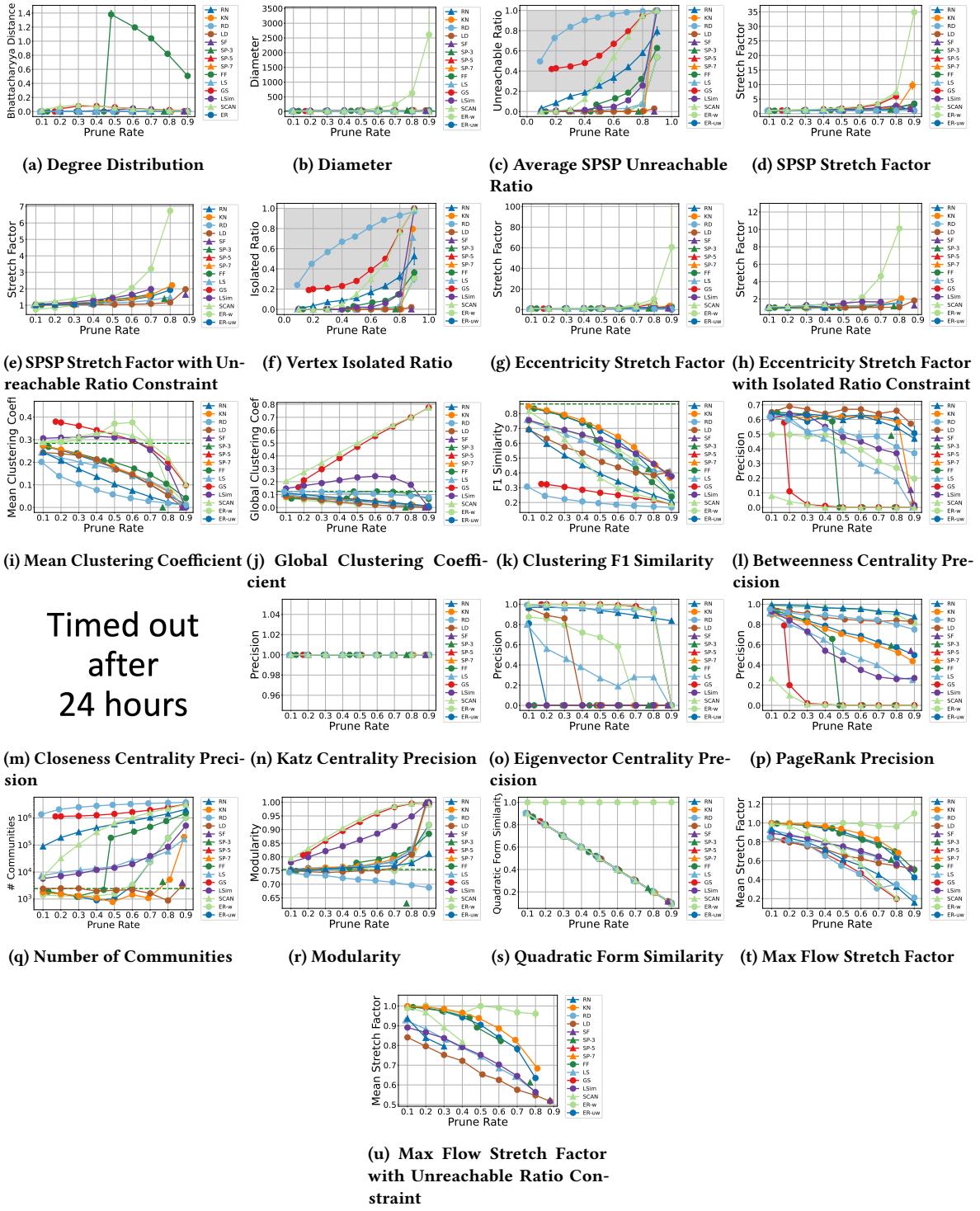


Figure 7: Metric Evaluation on com-LiveJournal

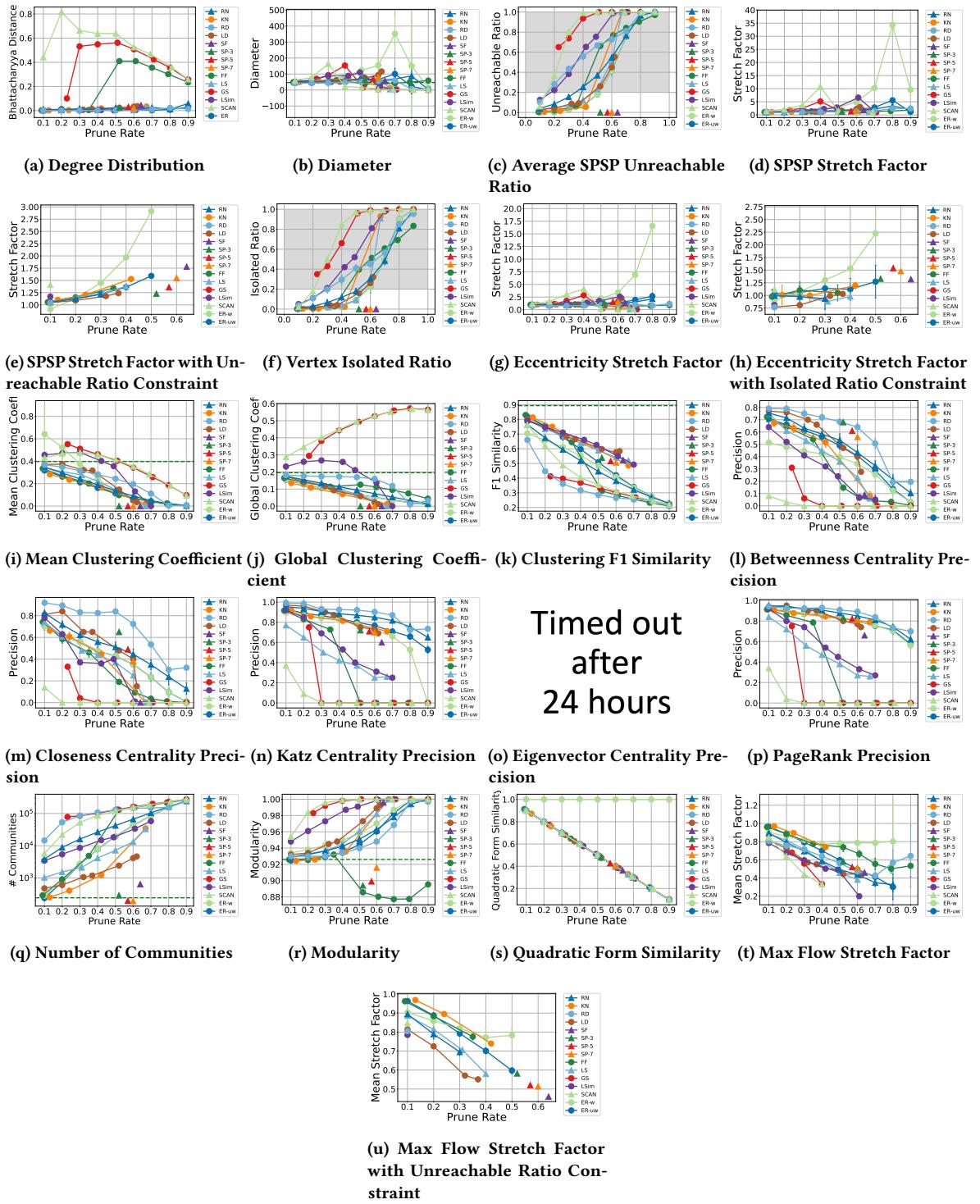


Figure 8: Metric Evaluation on com-Amazon

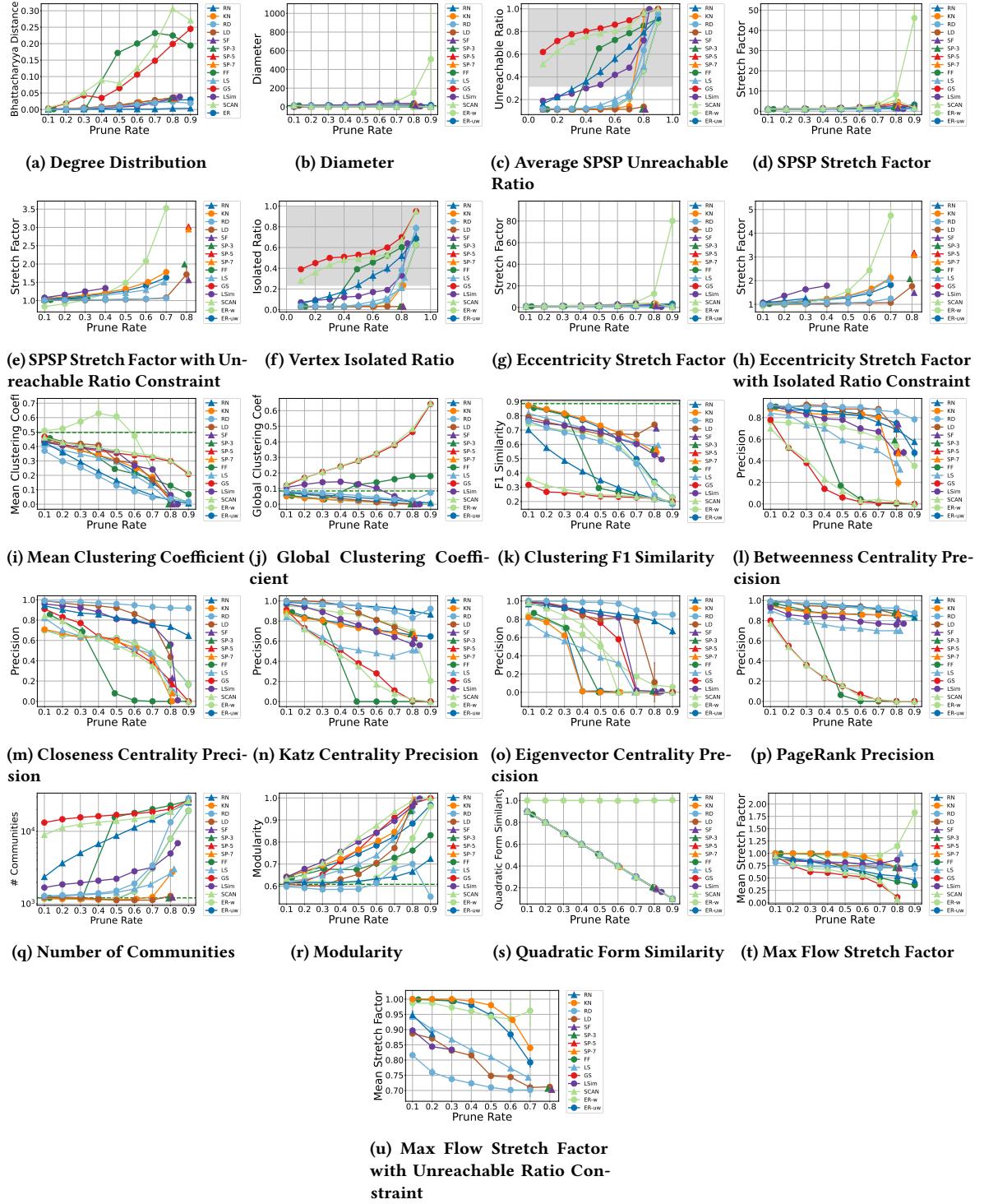


Figure 9: Metric Evaluation on email-Enron

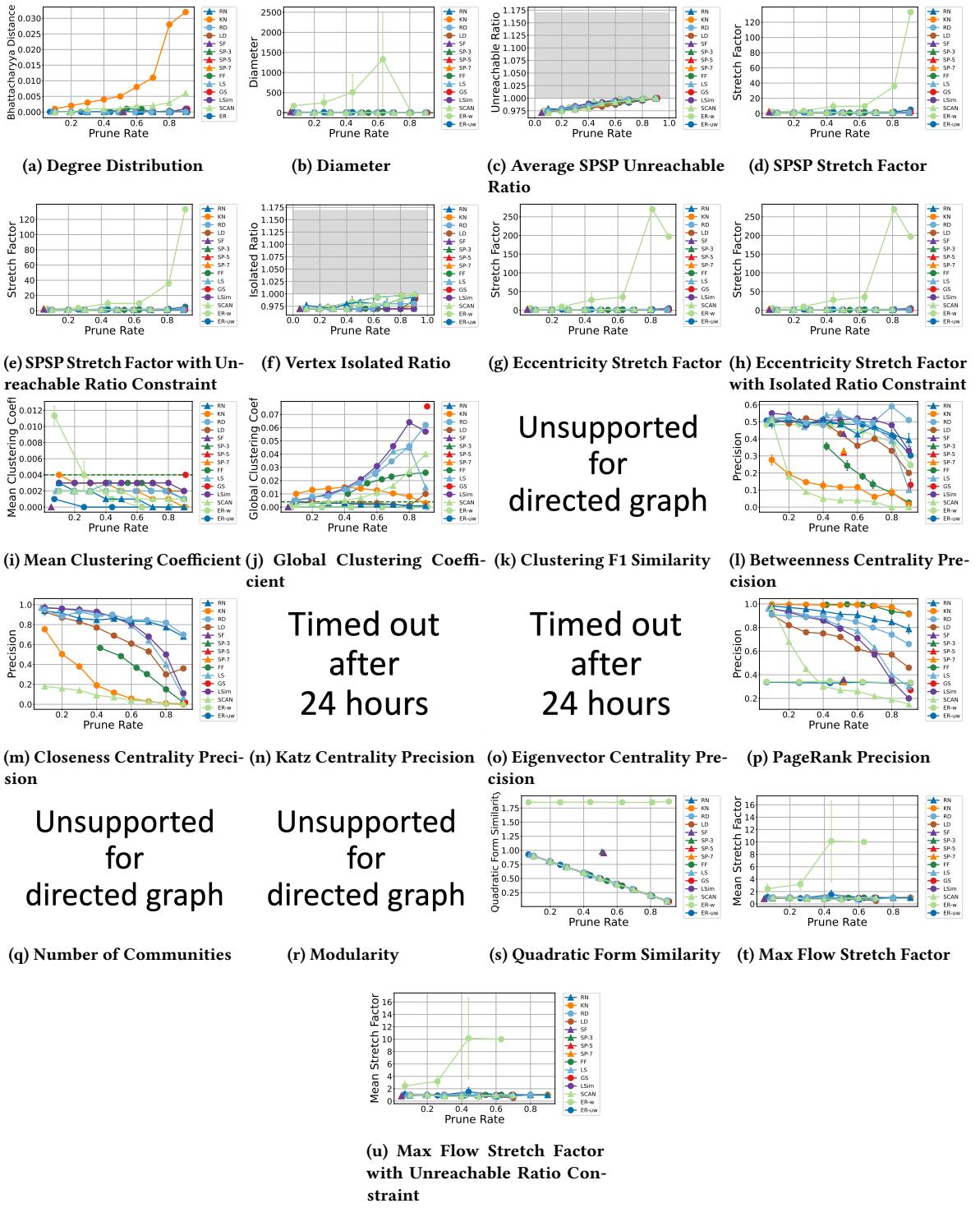


Figure 10: Metric Evaluation on wiki-Talk

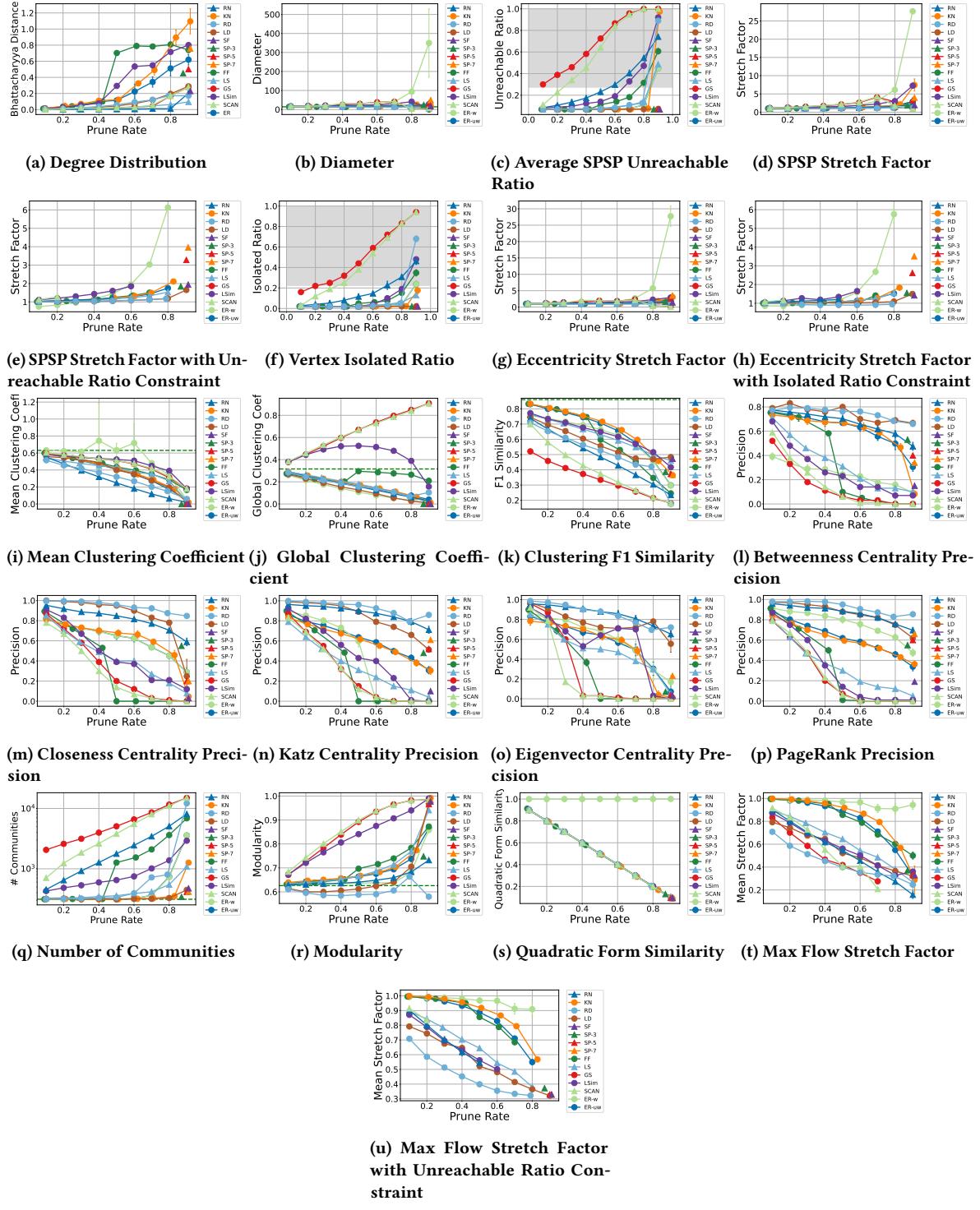


Figure 11: Metric Evaluation on ca-AstroPh

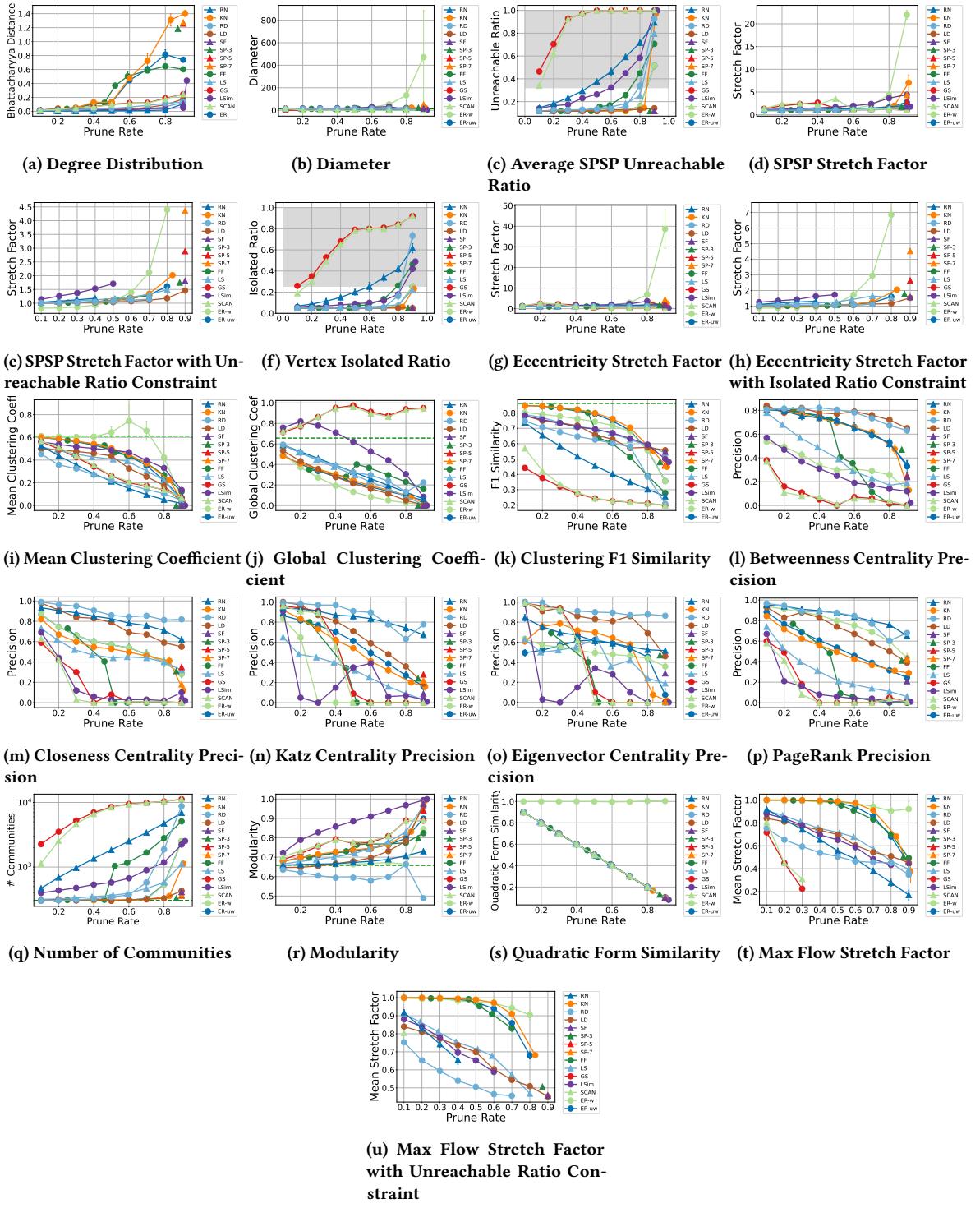


Figure 12: Metric Evaluation on ca-HepPh

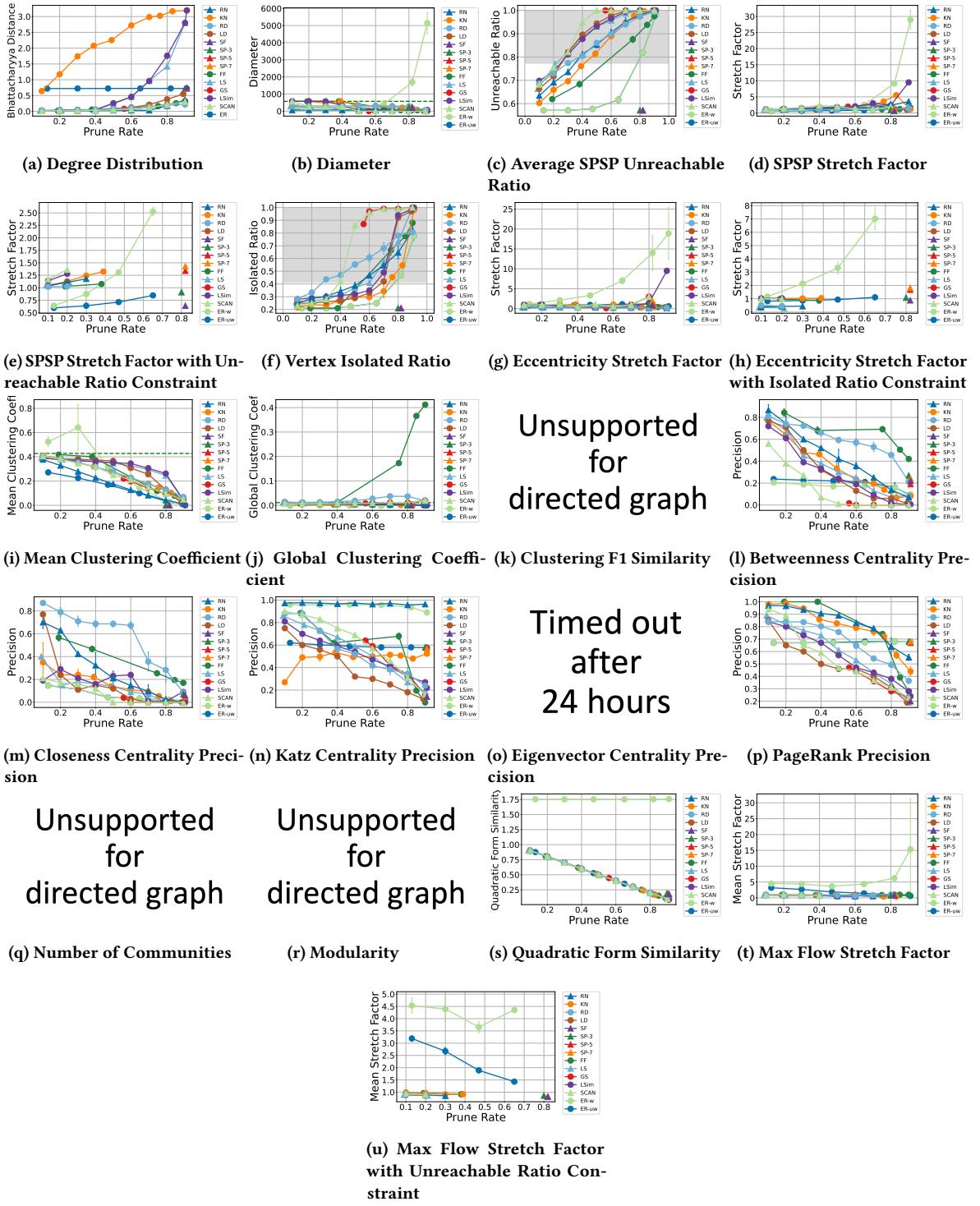


Figure 13: Metric Evaluation on web-BerkStan

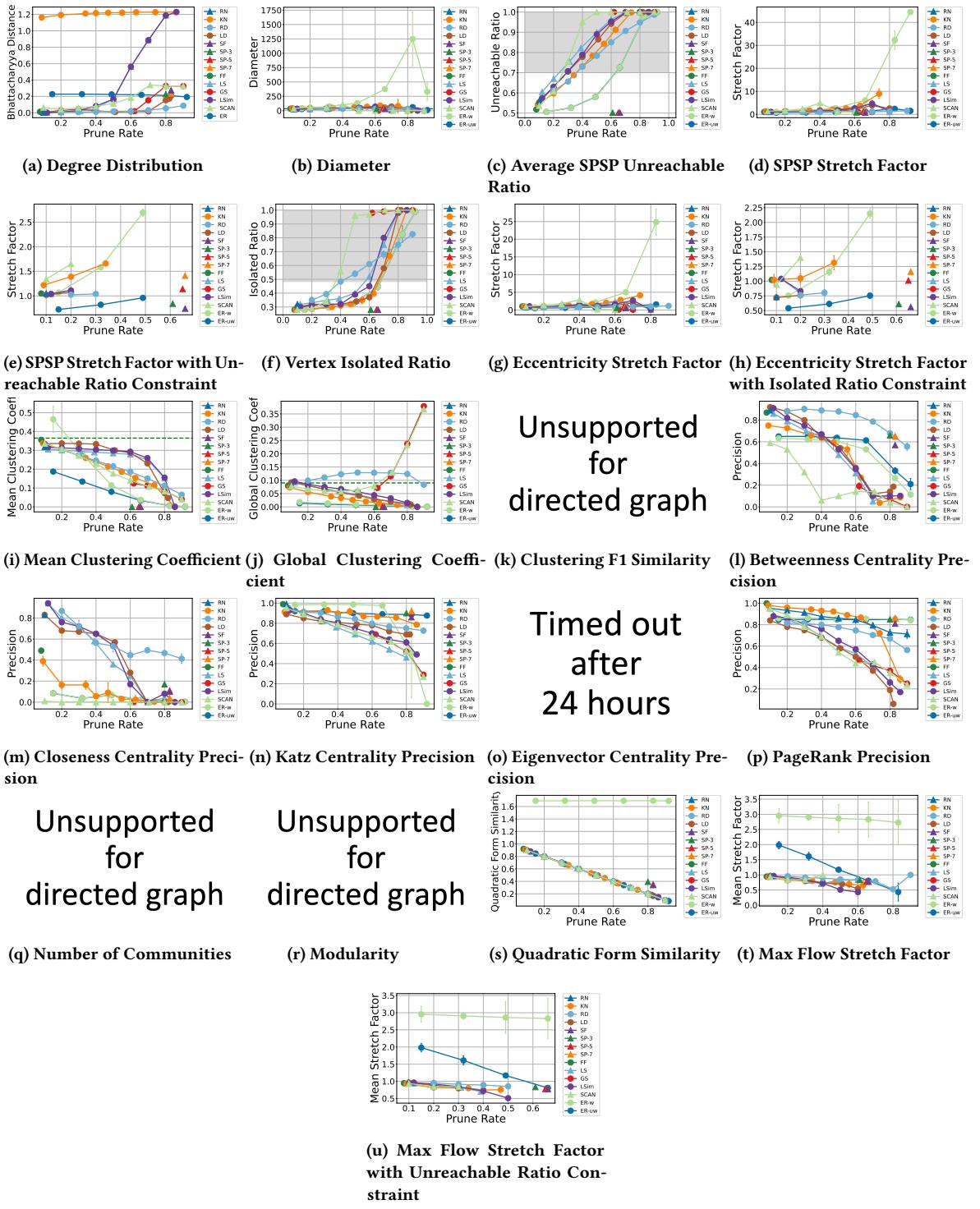


Figure 14: Metric Evaluation on web-Google

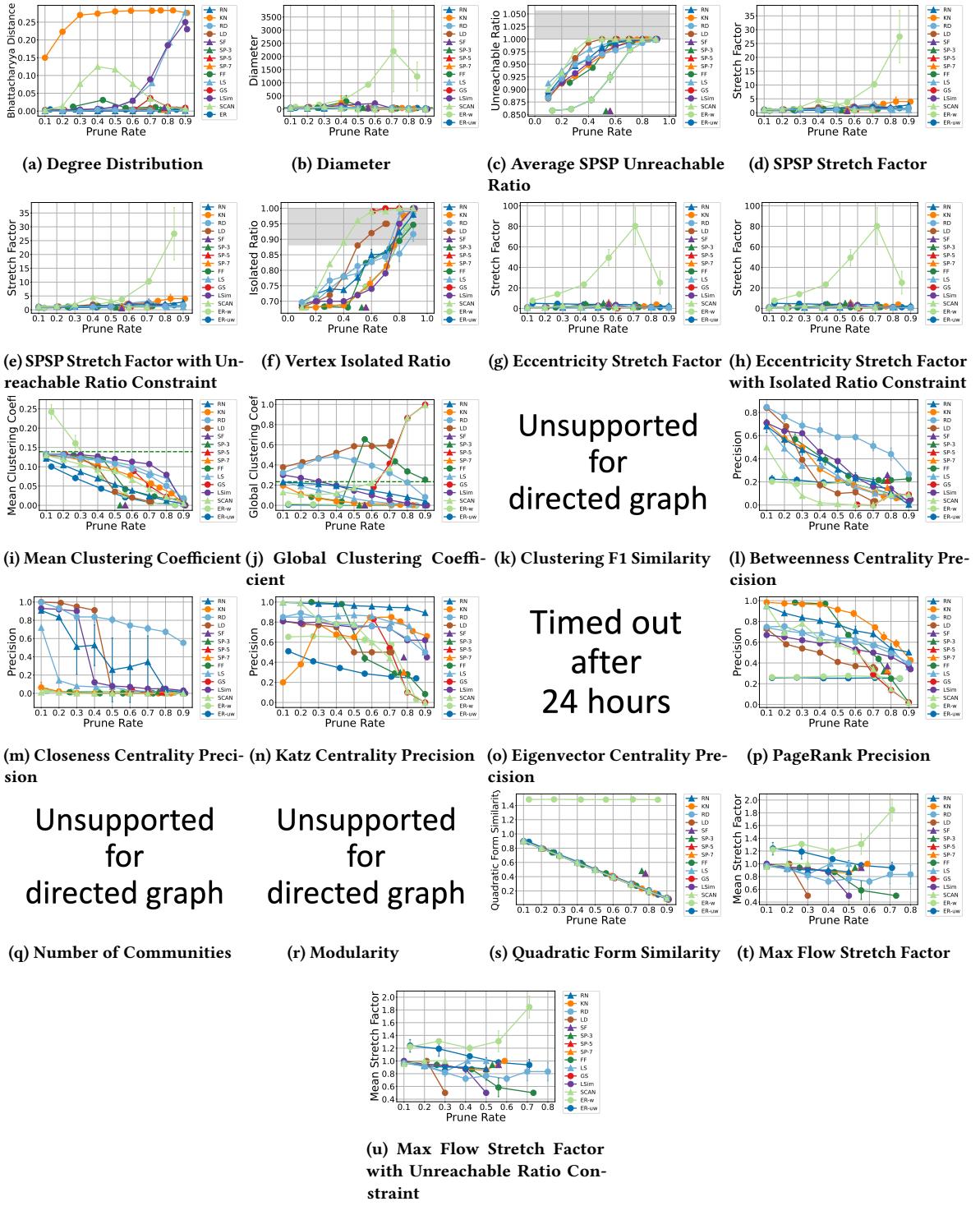


Figure 15: Metric Evaluation on web-NotreDame

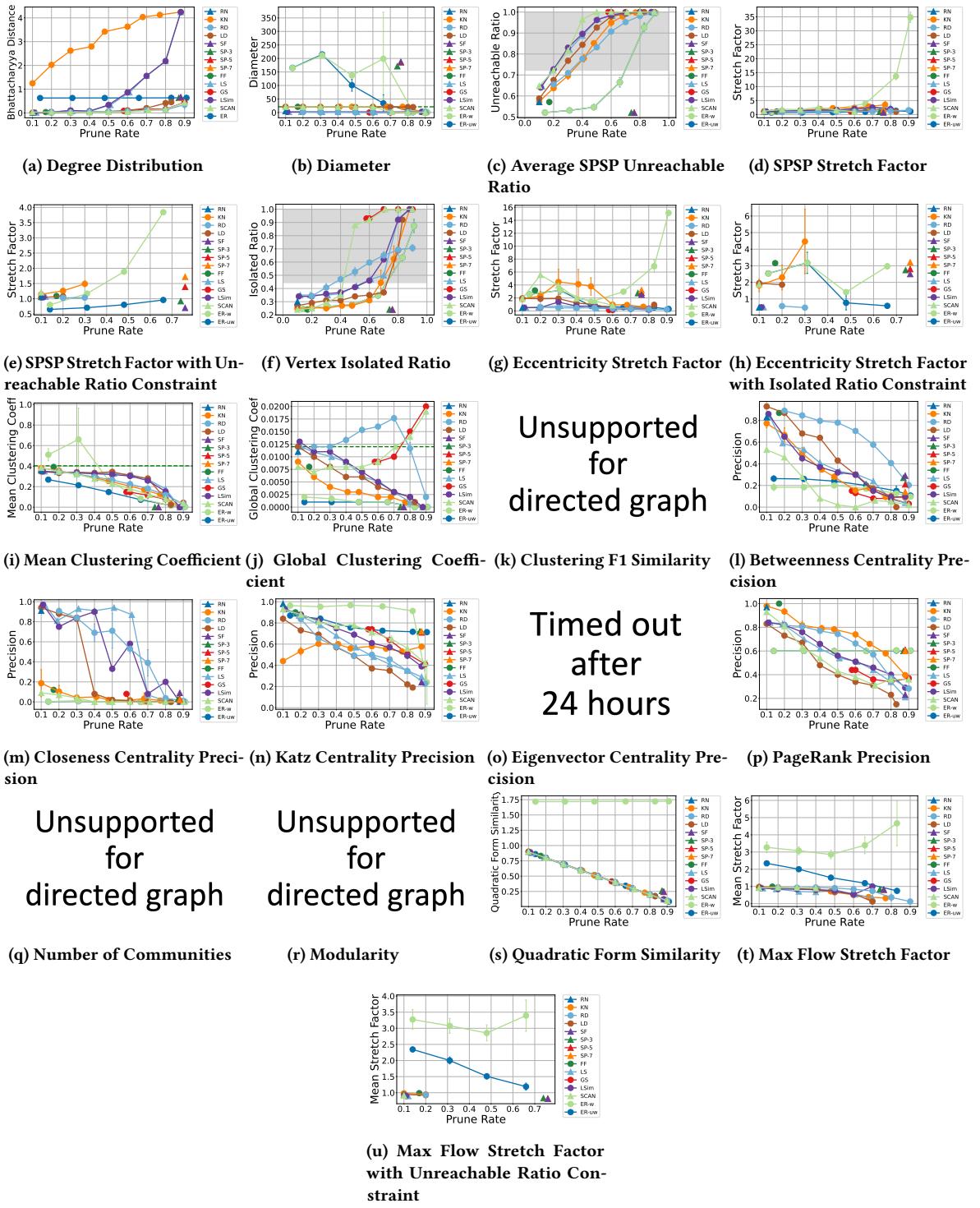


Figure 16: Metric Evaluation on web-Stanford

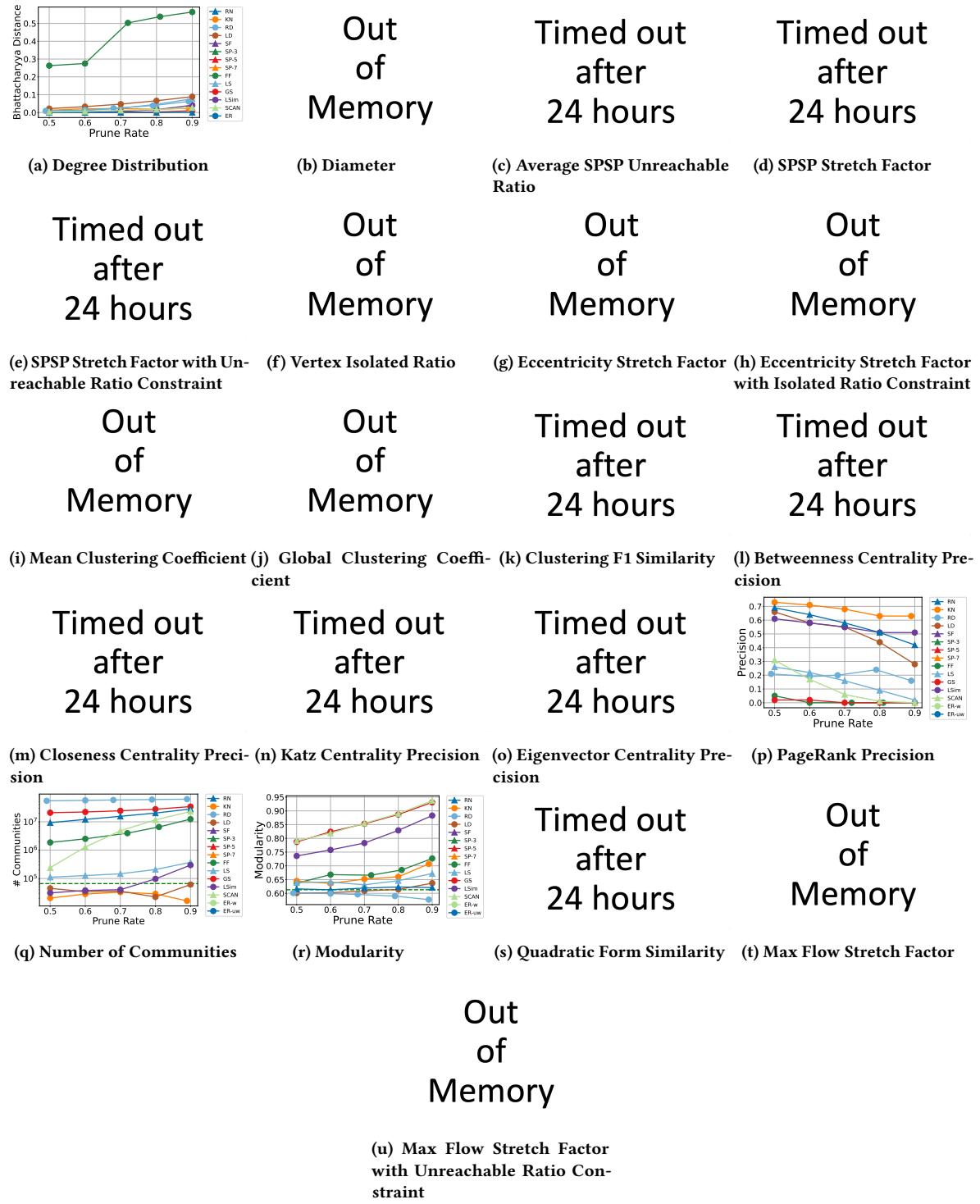


Figure 17: Metric Evaluation on com-friendster

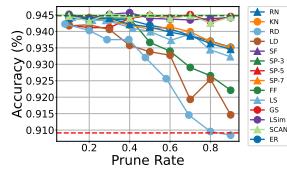


Figure 18: Clustering GCN Accuracy on Reddit

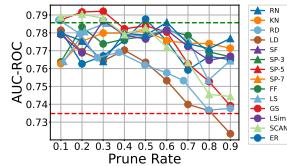


Figure 19: GraphSAGE Accuracy on ogbn-proteins