

Algorithms

May 24, 2016

Algorithm 1 Weight Assignment algorithm

Input:

1. Data frame $DF(from, to)$, each row signifying an edge between $from$ vertex and to vertex.

Output:

1. Vector wts containing the weight assigned to each edge in DF .

```
1: procedure WEIGHT( $DF(from, to)$ )
2:    $G(V, E) \leftarrow DF(from, to)$   $\triangleright E = (u, v) \forall u, v \in V$ 
3:    $simplify(G)$   $\triangleright simplify$  removes all self loops and multiple edges
4:    $deg \leftarrow degree(G)$ 
5:   Let  $wts$  be a vector of length  $|E|$  initialised to 0
6:   for each  $E_m = (u, v) \in E, m = 1, \dots, |E|$  do
7:      $start\_c \leftarrow deg[u]$ 
8:      $end\_c \leftarrow deg[v]$ 
9:      $wts[m] \leftarrow (start\_c + end\_c - 1) / |V| - 1$ 
10:  end for
11:  return  $wts$ 
12: end procedure
```

Algorithm 2 Link Prediction algorithm

Input:

1. Data frame $DF(from, to)$, each row signifying an edge between $from$ vertex and to vertex.
2. Floating-point value $threshold$ for the similarity score.

Output:

1. Data frame $sim(from, to, s)$, each row signifying an edge between $from$ vertex and to vertex with the corresponding similarity score s .

Ensure:

$G(V, E)$ is an incomplete graph.

```
1: procedure LINK( $DF(from, to), threshold$ )
2:    $G(V, E) \leftarrow DF(from, to)$   $\triangleright E = (u, v) \forall u, v \in V$ 
3:    $simplify(G)$   $\triangleright simplify$  removes all self loops and multiple edges
4:    $fc \leftarrow fastgreedy\_cluster(G)$   $\triangleright$ 
    $fastgreedy\_cluster(G)$  performs fastgreedy clustering algorithm on graph
    $G$  and returns each cluster as a subgraph  $C_i = (V_i, E_i), i = 1, \dots, c$  where  $c$ 
   is the number of identified clusters.
5:    $u \leftarrow K(V)$   $\triangleright K(V)$  returns the number of edges required for a
   complete graph of  $V$  vertices
6:    $le \leftarrow |E|/u$ 
7:    $lne \leftarrow u - |E|/u$ 
8:    $\Omega \leftarrow le/lne$ 
9:   for each  $u, v \in V$  do
10:    if  $u \neq v \wedge (u, v) \notin E$  then
11:       $total\_cn \leftarrow get\_common\_neighbours(u, v)$ 
12:       $u\_cluster \leftarrow fc\$membership[u]$   $\triangleright fc\$membership[u]$  returns
the cluster to which vertex  $u$  belongs to
13:       $v\_cluster \leftarrow fc\$membership[v]$ 
14:      if  $u\_cluster = v\_cluster$  then
15:        for each  $w \in V$  do
16:           $u\_cn < -get\_common\_neighbours(u, w)$ 
17:           $v\_cn < -get\_common\_neighbours(v, w)$ 
18:          if  $u\_cn \neq 0 \wedge v\_cn \neq 0$  then
19:            if  $fc\$membership[w] = u\_cluster$  then
20:               $within\_cn \leftarrow within\_cn + 1$ 
21:            end if
22:          end if
23:        end for
24:      else
25:         $within\_cn \leftarrow 0$ 
26:      end if
27:       $outside\_cn \leftarrow within\_cn - total\_cn$ 
28:      if  $within\_cn \neq 0 \wedge outside\_cn \neq 0$  then
29:         $s \leftarrow within\_cn/outside\_cn * \Omega$ 
30:      end if
31:      Let  $sim$  be a data frame having three columns, namely,  $from$ ,  $to$ 
and  $score$  initialised to 0
32:      if  $s > threshold$  then
33:         $append\_to\_data\_frame(sim, i, j, s)$ 
34:      end if
35:    end if
36:  end for
37:   $sort\_descending\_by\_column(sim, 3)$ 
38:  return  $sim$ 
39: end procedure
```

Algorithm 2 Link Prediction algorithm (continued)

```
40: procedure K( $G(V, E)$ )  
41:   Let  $num$  be an integer variable initialised to 0  
42:   if  $G$  is directed then  
43:      $num \leftarrow |V| * |V - 1|$   
44:   else  
45:      $num \leftarrow |V| * |V - 1|/2$   
46:   end if  
47:   return  $num$   
48: end procedure
```

Algorithm 3 Influence Analysis algorithm

Input:

1. Data frame $DF(from, to)$, each row signifying an edge between $from$ vertex and to vertex.
2. Integer $abscut_off$ which is the cutoff indegree for absolute cut score method.
3. Integer $fixcut_off$ which is the cutoff percentage for fixed percentage of population method.
4. Integer $sdcut_off$ which is the number of standard deviations for standard deviation method.
5. Integer num_sim which is the number of simulations to be performed in random permutation method.
6. Floating-point $randcut_off$ which is the p value for random permutation method.

Output:

1. Vector $abscut$ containing influentials identified by absolute cut score method.
2. Vector $fixcut$ containing influentials identified by fixed percentage of population method.
3. Vector $sdcut$ containing influentials identified by standard deviation method.
4. Vector $randcut$ containing influentials identified by random permutation method.

```
1: procedure INFLUENCE( $DF(from, to), abscut\_off, fixcut\_off, sdcut\_off, num\_sim, randcut\_off$ )
2:    $G(V, E) \leftarrow DF(from, to)$   $\triangleright E = (u, v) \forall u, v \in V$ 
3:    $simplify(G)$   $\triangleright simplify$  removes all self loops and multiple edges
4:    $C \leftarrow FIND\_STABLE\_CLUSTERED\_SAMPLE(G)$   $\triangleright$ 
    $FIND\_STABLE\_CLUSTERED\_SAMPLE(G)$  returns a set of clusters
   identified from a statistically significant sample
5:    $not\_yet \leftarrow G - C$ 
6:    $final \leftarrow CLUSTER\_BY\_PAGERANK(C, not\_yet)$   $\triangleright$ 
    $CLUSTER\_BY\_PAGERANK(C, not\_yet)$  returns a set of clusters which
   have all vertices of  $not\_yet$  clustered to one of the clusters in  $C$  based on
   pagerank algorithm
7:   Let  $abscut, fixcut, sdcut, randcut$  be vectors initialised to 0.
8:   for each cluster( $V_{cluster}, E_{cluster}$ )  $\in final$  do
9:      $in\_ties \leftarrow in\_degree(cluster)$ 
10:     $desc\_in\_ties \leftarrow sort\_decreasing(in\_ties)$ 
11:     $scree\_plot \leftarrow plot(desc\_in\_ties)$ 
12:     $draw\_horizontal\_abline(scree\_plot, abscut\_off)$ 
13:     $append\_points\_above\_abline\_to\_vector(abscut)$ 
14:     $cut\_off\_fix \leftarrow fixcut\_off * |V_{cluster}|$ 
15:     $draw\_vertical\_abline(scree\_plot, cut\_off\_fix)$ 
16:     $append\_points\_left\_of\_abline\_to\_vector(fixcut)$ 
17:     $cut\_off\_sd \leftarrow mean(in\_ties) + sdcut\_off * sd(in\_ties)$ 
18:     $draw\_horizontal\_abline(scree\_plot, cut\_off\_sd)$ 
19:     $append\_points\_above\_abline\_to\_vector(sdcut)$ 
20:     $sim \leftarrow simulation\_conditional\_on\_outdegree(cluster, num\_sim)$   $\triangleright$ 
     $simulation\_conditional\_on\_outdegree(cluster, num\_sim)$  performs random
    simulation of  $cluster$  graph  $num\_sim$  number of times constraining on the
    outdegree.
21:    Let  $sim\_in\_ties$  be a vector of length  $|V_{cluster}| * num\_sim$ 
22:    for  $j \leftarrow 1$  to  $\#\{sim\}$  do  $\triangleright \#\{sim\}$  returns the number of clusters
    in  $sim$ 
23:       $append\_to\_vector(sim\_in\_ties, indegree(sim[j]))$ 
24:    end for
```

Algorithm 3 Influence Analysis algorithm (continued)

```
25:    $q \leftarrow \text{quantile}(\text{sim\_in\_ties}, \text{randcut\_off})$   $\triangleright$   
    $\text{quantile}(\text{sim\_in\_ties}, \text{randcut\_off})$  returns sample quantiles of  $\text{sim\_in\_ties}$   
   corresponding to the given probability  $\text{randcut\_off}$   
26:    $\text{draw\_histogram}(\text{sim\_in\_ties})$   
27:    $\text{draw\_vertical\_abline}(\text{sim\_in\_ties}, q)$   
28:    $\text{append\_points\_right\_of\_abline\_to\_vector}(\text{randcut})$   
29: end for  
30: return  $\text{abscut}, \text{fixcut}, \text{sdcut}, \text{randcut}$   
31: end procedure  
  
32: procedure FIND_STABLE_CLUSTERED_SAMPLE( $G(V, E)$ )  
33:    $\text{perc} \leftarrow 0.1$   
34:   Let  $\text{major\_fc\_length}$ ,  $\text{temp}$  and  $\text{diff}$  be vectors of length 10 initialised  
   to 0.  
35:    $i \leftarrow 1$   
36:   while  $\text{perc} \leq 1$  do  
37:      $\text{len} \leftarrow |V| * \text{perc}$   
38:      $\text{sample} \leftarrow \text{random\_sample\_without\_repetition}(G, \text{len})$   $\triangleright$   
      $\text{random\_sample\_without\_repetition}(G, \text{len})$  randomly selects  $\text{len}$  number of  
     vertices and edges corresponding to these vertices from graph  $G$   
39:      $\text{fc} \leftarrow \text{fastgreedy\_cluster}(\text{sample})$   $\triangleright \text{fastgreedy\_cluster}(\text{sample})$   
     performs fastgreedy clustering algorithm on graph  $\text{sample}$  and returns each  
     cluster as a subgraph  $C_i = (V_i, E_i), i = 1, \dots, c$  where  $c$  is the number of  
     identified clusters.  
40:      $\text{major\_fc} \leftarrow \text{fc}[\text{sizes}(\text{fc}) > (0.01 * \text{len})]$   $\triangleright \text{sizes}(\text{fc})$  returns the  
     number of vertices in each cluster present in  $\text{fc}$   
41:      $\text{major\_fc\_length}[i] \leftarrow \#\{\text{major\_fc}\}$   
42:      $\text{temp}[i] \leftarrow \text{major\_fc\_length}[i]$   
43:      $\text{res} \leftarrow 0$   
44:     if  $i > 1$  then  
45:        $\text{diff}[i - 1] \leftarrow |\text{temp}[i] - \text{temp}[i - 1]|$   
46:       if  $i > 2$  then  
47:         if  $|\text{diff}[i - 1] - \text{diff}[i - 2]| \leq \text{res}$  then  
48:            $\text{res} \leftarrow |\text{diff}[i - 1] - \text{diff}[i - 2]|$   
49:           break  
50:         else  
51:            $\text{res} \leftarrow |\text{diff}[i - 1] - \text{diff}[i - 2]|$   
52:         end if  
53:       end if  
54:     end if  
55:      $i++$   
56:      $\text{perc} \leftarrow \text{perc} + 0.1$   
57:   end while  
58:   return  $\text{major\_fc}$   
59: end procedure
```

Algorithm 3 Influence Analysis algorithm (continued)

```
60: procedure CLUSTER_BY_PAGERANK( $C(V_c, E_c), not\_yet(V'_c, E'_c)$ )
61:   Let pageranks be vector of length  $\#\{C\}$  initialised to 0.
62:   while  $\#|V_c| \neq 0$  do
63:      $random(key, E_{key}) \leftarrow randomly\_select\_subgraph\_without\_repetition(not\_yet, 1)$ 
      $\triangleright randomly\_select\_subgraph\_without\_repetition(not\_yet, 1)$  randomly se-
     lects 1 vertex, key from not_yet and corresponding edges  $E_{key} = (key, V_o)$ ,
      $V_o$  being vertices in not_yet having edge to key.
64:      $delete\_from\_graph(not\_yet, random)$ 
65:     for each  $C_i = (V_i, E_i) \in C, i = 1, \dots, \#\{C\}$  do
66:       Let  $check(V_{check}, E_{check})$  be an empty graph object.
67:        $check \leftarrow C_i$ 
68:        $append\_to\_list(V_{check}, key)$ 
69:       for each  $E_{jkey} = (key, V_{jo}) \in E_{key}, j = 1, \dots, |E_{key}|$  do
70:         if  $V_{jo} \in V_{check}$  then
71:            $append\_to\_list(E_{check}, E_{jkey})$ 
72:         end if
73:       end for
74:        $pr \leftarrow pagerank(check)[key]$ 
75:        $pageranks[i] \leftarrow pr$ 
76:     end for
77:     if  $sum(pageranks) > 0$  then
78:        $ind \leftarrow index(max(pageranks))$ 
79:        $append\_to\_list(V_{ind}, key)$ 
80:       for each  $E_{jkey} = (key, V_{jo}) \in E_{key}, j = 1, \dots, |E_{key}|$  do
81:         if  $V_{jo} \in V_{ind}$  then
82:            $append\_to\_list(E_{ind}, E_{jkey})$ 
83:         end if
84:       end for
85:     end if
86:   end while
87:   return  $C$ 
88: end procedure
```

Algorithm 4 Time Series Analysis algorithm

Input:

1. Data frame $DF(from, to, time)$, each row signifying an edge between $from$ vertex and to vertex with the corresponding time, $time$
2. Floating-point value $start$ indicating the value of $time$ when analysis must commence.
3. Floating-point value end indicating the value of $time$ when analysis must terminate.
4. Floating-point value $increment$ indicating the increments of $start$ in every iteration.

Output:

1. Set of plots, $images(plots)$ containing the plots of the graph at every iteration.

```
1: procedure TIME( $DF(from, to, time), start, end, increment$ )
2:    $G(V, E, T) \leftarrow DF(from, to, time)$   $\triangleright E = (u, v, t) \forall u, v \in V, \forall t \in T$ 
3:    $vcolor \leftarrow generate\_color\_palette(V)$   $\triangleright generate\_color\_palette(V)$ 
   returns a distinct color for each vertex in  $V$ 
4:    $i \leftarrow start$ 
5:   while  $i \leq end$  do
6:     for each  $e \in E, v \in V$  do
7:       if  $e$time < i$  then
8:          $e$weight \leftarrow 1$ 
9:          $e$color \leftarrow "gray"$ 
10:      else
11:         $e$weight \leftarrow 0$ 
12:         $e$color \leftarrow "black"$ 
13:      end if
14:      if  $strength(v) = 0$  then  $\triangleright strength$  returns the sum of edge
        weights of adjacent edges of vertex  $v$ 
15:         $v$color \leftarrow "black"$ 
16:      else
17:         $v$color \leftarrow vcolor[v]$ 
18:      end if
19:       $v$size \leftarrow 1 + 2 * \log(strength(v))$ 
20:    end for
21:     $append\_to\_plots(images, plots)$ 
22:     $i \leftarrow i + increment$ 
23:  end while
24: end procedure
```
