

Clustering Via Connections

by
Jonathan McFadden

6/13/2016 -

Table of Contents

<u>1. Definitions</u>	1-	6/14
• 1.1 - Input Matrix	1-2	6/14
• 1.2 - SVD Matrices	2-4	6/14
• 1.3 - Connection Matrices	4-8	6/14
• 1.4 - Relation sets	8-12	6/14
• 1.5 - Cluster sets	12-15	6/14

1 Definitions

1.1 - Input Matrix

Let $\bar{L} \in \mathbb{R}^{m \times n}$ be a matrix that represents how m terms describe a collection of n objects. We express this formally as either

$$L_{ij} = \begin{cases} 1 & \text{if the } i^{\text{th}} \text{ term describes the } j^{\text{th}} \text{ object.} \\ 0 & \text{otherwise} \end{cases}$$

or

$$L_{ij} = \begin{cases} n & \text{if the } i^{\text{th}} \text{ terms describes the } j^{\text{th}} \text{ object } n \text{ times} \\ 0 & \text{otherwise} \end{cases}$$

In both definitions for the L_{ij}
we require that $i \in [1, m] \subset \mathbb{Z}^+$
and that $j \in [1, n] \subset \mathbb{Z}^+$.

1.2 - SVD Matrices

we will be using both the
full SVD of \bar{L} and the k
order reduced SVD of \bar{L} . The
order of the reduced SVD, k ,
must satisfy

$$k < r$$

where $r = \text{RANK}[\bar{L}]$.

1.2.1 - Full SVD

Let the full singular value decomposition of \bar{L} be given by

$$\bar{L} = \bar{U} \bar{S} \bar{V}^T$$

where $\bar{U} \in \mathbb{R}^{m \times r}$, $\bar{S} \in \mathbb{R}^{r \times r}$, and $\bar{V} \in \mathbb{R}^{n \times r}$, with r defined as above. We call \bar{U} the 'row matrix' of \bar{L} , \bar{S} the 'singular value matrix' of \bar{L} , and \bar{V} the 'column matrix' of \bar{L} .

1.2.2 - Reduced SVD

Let the K order reduced singular

value decomposition of \bar{L} be given by

$$\bar{L} \cong \bar{U}_x \bar{S}_x \bar{V}_x^T$$

where $\bar{U}_x \in \mathbb{R}^{m \times X}$, $\bar{S}_x \in \mathbb{R}^{X \times X}$, and $\bar{V}_x \in \mathbb{R}^{n \times X}$ with X defined as above. We call \bar{U}_x the 'reduced row matrix' of \bar{L} , \bar{S}_x the 'reduced singular value matrix' of \bar{L} , and \bar{V}_x the 'reduced column matrix' of \bar{L} .

1.3 - Connection Matrices

From \bar{L} , we will create two

5/

'Connection Matrices'. These matrices will be used to provide information for clustering either the rows of \bar{E} or the columns of \bar{E} .

1.3.1 - Row Connection Matrix

The first of the 'Connection Matrices' is the 'Row Connection Matrix'. This matrix will be used to provide information for clustering the rows of \bar{E} . We will denote this matrix by $\bar{A}_v \in \mathbb{R}^{m \times m}$ and define it, in terms of \bar{E} , by either

$$\bar{\bar{A}}_v = \bar{\bar{L}} \bar{\bar{L}}^+$$

or

$$(\bar{\bar{A}}_v)_{ij} = \sum_{k=1}^n \left\{ S^*(\bar{\bar{L}}, i, k) \bar{\bar{L}}_{jk} \right\}$$

where $i, j \in [1, m] \subset \mathbb{Z}^+$ and the function $S^*(\bar{\bar{L}}, l, \lambda)$ is defined

$$S^*(\bar{\bar{L}}, l, \lambda) = \begin{cases} 1 & \text{if } \bar{\bar{L}}_{l\lambda} \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

with $l \in [1, m] \subset \mathbb{Z}^+$ and $\lambda \in [1, n] \subset \mathbb{Z}^+$.
for the definition of $S^*(\bar{\bar{L}}, l, \lambda)$.

1.3.2 - Column Connection Matrix

The second 'Connection Matrix' is the 'Column Connection Matrix'. This matrix will provide information when clustering the columns of \bar{L} and is denoted $\bar{A}_v \in \mathbb{R}^{n \times n}$. We define this matrix, in terms of \bar{L} , by either

$$\bar{A}_v = \bar{L}^T \bar{L}$$

or

$$(\bar{A}_v)_{ij} = \sum_{k=1}^m \left\{ S^*(\bar{L}, k, i) \bar{L}_{kj} \right\}$$

where $i, j \in [1, n] \subset \mathbb{Z}^+$ and $S^*(\bar{z}, z)$ is as defined above.

1.4 - Relation Sets

We will now define two sets of 2-tuples, with one for rows and the other for columns. These sets will describe the relation or relations any row/column element with any other row/column element as the strength of the relation or relations in question.

1.4.1 - Row Relation Set

Let \mathbb{U} be a set of m 2-tuples such that

$$\mathbb{U} = \left\{ \{M_i, U_i\}, \{M_2, U_2\}, \dots, \{M_m, U_m\} \right\}$$

where the $\{M_i, U_i\}$ are 2-tuples representing all of the rows to which the i^{th} row is related, U_i , and the strength of those relations, M_i , with $i \in [1, m] \subset \mathbb{Z}^+$. Let us now define U_i and M_i concurrently as

$$U_i = \left\{ j \in [1, m] \subset \mathbb{Z}^+ \mid (\bar{A}_U)_{ij} \neq \emptyset \right\}$$

and

$$\mathbb{M}_i = \{(\bar{A}_o)_{ij} \mid j \in [1, m] \subset \mathbb{Z}^+ \text{ and } (\bar{A}_o)_{ij} \neq \emptyset\}$$

These definitions require that

$$|\mathbb{M}_i| = |\mathbb{U}_i|$$

for all $i \in [1, m] \subset \mathbb{Z}^+$. Additionally, the strength of the connection to the l^{th} element in \mathbb{U}_i is represented by the l^{th} element in \mathbb{M}_i .

1.4.2 - Column Relation Set

Let \mathbb{V} be a set of n 2-tuples such that

$$\mathbb{V} = \{\{N_1, V_1\}, \{N_2, V_2\}, \dots, \{N_n, V_n\}\}$$

where the $\{N_j, V_j\}$ are 2-tuples representing all of the columns to which the j^{th} column is related, V_j , and the strength of those relations, N_j ; with $j \in [1, n] \subset \mathbb{Z}^+$. Let us now define V_j and N_j concurrently as

$$V_j = \{i \in [1, n] \subset \mathbb{Z}^+ \mid (\bar{A}_v)_{ji} \neq 0\}$$

and

$$N_j = \{(\bar{A}_v)_{ji} \mid i \in [1, n] \subset \mathbb{Z}^+ \text{ and } (\bar{A}_v)_{ji} \neq 0\}$$

These definitions require that

$$|\mathbb{N}_j| = |\mathbb{V}_j|$$

for all $j \in [1, n] \subset \mathbb{Z}^+$. Additionally, the strength of the connection to the l^{th} element in \mathbb{V}_j is represented by the l^{th} element in \mathbb{N}_j .

1.5 - Cluster Sets

We now define two sets of sets, with one set representing the row clusters and the other representing the column clusters.

These sets consist of sets that represent the members of each cluster.

1.5.1 - Row Cluster Set

The 'Row Cluster Set' is a set of u_m sets representing each of the u_m clusters. We denote this \mathbb{U}^* and describe it by

$$\mathbb{U}^* = \{ U_1^*, U_2^*, \dots, U_{u_m}^* \}$$

where the U_u^* are sets containing the members of the u^{th} cluster, with $u \in [1, u_m] \subset \mathbb{Z}^+$. Each U_u^* is constructed from the $\{\mathbb{M}_i, \mathbb{W}_i\}$ of

its (the V_n^*) constituent elements,
via a process which is
defined later.

1.5.2 - Column Cluster Set

The 'Column Cluster Set' is a set of v_n sets representing each of the v_n clusters. We will denote this V^* and describe it by

$$V^* = \{V_1^*, V_2^*, \dots, V_{v_n}^*\}$$

where the V_n^* are sets containing the members of the n^{th} cluster, with $n \in [1, v_n] \subset \mathbb{Z}^+$. Each of the V_n is constructed from

the $\{N_j, V_j\}$ of its (the V_j^*) constituent elements through a process which will be described later.

2 Disjoint Sets

When \bar{L} consists of a series of disjoint sets, the algorithm proceeds by following the steps

- 1) Compute \bar{A}_o
- 2) Construct $U = \{\{M_1, U_1\}, \dots, \{M_m, U_m\}\}$
- 3) Complete the U_i
- 4) Construct the U^* set
- 5) Compute \bar{A}_v
- 6) Construct $V = \{\{N_1, V_1\}, \dots, \{N_n, V_n\}\}$
- 7) Complete the V_j
- 8) Construct the V^* set
- 9) Reorder \bar{L} into a new matrix, \bar{L}^*

2.1 - Compute \bar{A}_{uv}

Using the expression

$$\bar{A}_{uv} = \bar{L} \bar{L}^+$$

compute $\bar{A}_{uv} \in \mathbb{R}^{m \times m}$. Alternatively,
the expression

$$(\bar{A}_{uv})_{ij} = \sum_{k=1}^n \{ S^*(\bar{L}, i, k) (\bar{L})_{jk} \}$$

may be used if different
weighting is desired.

2.2 - Construct V

We initially construct V by

the following routine

```

FOR ( $i = 1:m$ ) {
     $U_i = \emptyset$ ;  $M_i = \emptyset$ ;
    FOR ( $j = 1:m$ ) {
        IF ( $(\bar{A}_v)_{ij} \neq 0$ ) {
             $U_i = U_i \cup \{j\}$ ;
             $M_i = M_i \cup \{(\bar{A}_v)_{ij}\}$ ;
        }
    }
     $U[i, 1] = M_i$ ;  $U[i, 2] = U_i$ ;
}

```

After completion of this routine, we make a copy of this initial U , denoted U_0 .

2.3 - Complete the U sets

The elements of \mathbb{U} must now be 'completed' to include indirectly related elements. Taking advantage of the disjointedness of \mathbb{L} , this 'completion' is accomplished using the simple routine below

BOOLEAN ISCHANGED = TRUE;

WHILE (ISCHANGED) {
 ISCHANGED = FALSE;

FOR ($i = 1 : m$) {

| $U_i = U[i, 2]$;

| $M_i = U[i, 1]$;

| FOR ($j = 1 : m$) {

| | $U_j = U[j, 2]$;

| | $M_j = U[j, 1]$;

```

    |   IF ( $U_i \neq U_j$  AND  $U_i \cap U_j \neq \emptyset$ ) {
    |       :  $U_i = U_i \cup \{U_j \setminus U_i\};$ 
    |       :  $M_i = M_i \cup \{M_j \setminus M_i\};$ 
    |       :  $U_j = U_j \cup \{U_i \setminus U_j\};$ 
    |       :  $M_j = M_j \cup \{M_i \setminus M_j\};$ 
    |       : ISCHANGED = TRUE;
    |   }
    |   : END
    |
    |   IF (ISCHANGED) {
    |       :  $U[j, 2] = U_j;$ 
    |       :  $U[j, 1] = M_j;$ 
    |   }
    |   : END
    |
    |   IF (ISCHANGED) {
    |       :  $U[i, 2] = U_i;$ 
    |       :  $U[i, 1] = M_i;$ 
    |   }
    |   : END
    |
    : END

```

2.4 - Construct the U^* sets

To construct the U^* sets we look for unique U_i in U and store each unique U_i as its own set in U^* . This is done via the following routine

```
U*[1] = U[1,2]; CLUSTCNT = 1;  
FOR (i=2:m) {  
    Ui = U[i,2]; BOOLEAN ISNEW = TRUE;  
  
    FOR (j = 1: CLUSTCNT) {  
        U*j = U*[j];  
  
        IF (U*j == Ui) {  
            ISNEW = FALSE;  
            BREAK;  
        } END  
    } END  
  
    IF (ISNEW) {  
        CLUSTCNT = CLUSTCNT + 1;  
        U*[CLUSTCNT] = Ui;  
    } END  
} END
```

2.5 - Compute \bar{A}_v

Using the expression

$$\bar{A}_v = \bar{L}^T \bar{L}$$

compute $\bar{A}_v \in \mathbb{R}^{n \times n}$. Alternatively,
the expression

$$(\bar{A}_v)_{ij} = \sum_{k=1}^m \{ S^*(\bar{L}_{ik}) (\bar{L})_{kj} \}$$

may be used if different
weighting is desired.

2.6 - Construct V

We initially construct V by the

following routine

```

FOR ( $i = 1:n$ ) {
     $V_i = \emptyset$ ;  $N_i = \emptyset$ ;
    FOR ( $j = 1:n$ ) {
        IF  $((\bar{A}_v)_{ij} \neq 0$ ) {
             $V_i = V_i \cup \{j\}$ ;
             $N_i = N_i \cup \{(\bar{A}_v)_{ij}\}$ ;
        }
    }
     $\} END$ 
}
 $V[i, 1] = N_i$ ;  $V[i, 2] = V_i$ ;
 $} END$ 

```

After the completion of this routine, we make a copy of this initial V , denoted V_0 .

2.7 - Complete the ∇ Sets

The elements of ∇ must now be 'completed' to include indirectly related elements.

Taking advantage of \bar{L} 's disjointedness, this 'completion' can be done using the simple routine below

BOOLEAN ISCHANGED = TRUE;

WHILE (ISCHANGED) {

 ISCHANGED = FALSE;

 FOR ($i = 1 : n$) {

$V_i = \nabla[i, 2]$;

$DV_i = \nabla[i, 1]$;

 FOR ($j = 1 : n$) {

$V_j = \nabla[j, 2]$;

$DV_j = \nabla[j, 1]$;

```

|   |   | IF ( $V_i \neq V_j$  AND  $V_i \cap V_j \neq \emptyset$ ) {
|   |   |   |   |  $V_i = V_i \cup \{V_j \setminus V_i\};$ 
|   |   |   |   |  $N_i = N_i \cup \{N_j \setminus N_i\};$ 
|   |   |   |   |  $V_j = V_j \cup \{V_i \setminus V_j\};$ 
|   |   |   |   |  $N_j = N_j \cup \{N_i \setminus N_j\};$ 
|   |   |   |   | ISCHANGED = TRUE;
|   |   | } END
|   |   | IF (ISCHANGED) {
|   |   |   |   |  $V[j, 2] = V_j;$ 
|   |   |   |   |  $V[j, 1] = N_j;$ 
|   |   | } END
|   |   | IF (ISCHANGED) {
|   |   |   |   |  $V[i, 2] = V_i;$ 
|   |   |   |   |  $V[i, 1] = N_i;$ 
|   |   | } END
|   | } END
| } END

```

2.8 - Construct the V^* sets

To construct the sets in V^* , we look for unique V_j in V , then store each unique V_j as its own set in V^* . This is done via the following routine

```
V^*[1] = V[1, 2]; CLUSTCNT = 1;  
FOR (i = 2:n) {  
    BOOLEAN ISNEW = TRUE;  
    V_i = V[i, 2];  
  
    FOR (j = 1:CLUSTCNT) {  
        V^*_j = V^*[j];  
        IF (V^*_j == V_i) {  
            ISNEW = FALSE; BREAK;  
        }  
    }  
    IF (ISNEW) {  
        CLUSTCNT++; V^*[CLUSTCNT] = V_i;  
    }  
}
```

2.9 - Reorder \bar{L} into \bar{L}^*

The last step is to reorder \bar{L} into the new, clustered matrix $\bar{L}^* \in \mathbb{R}^{m \times n}$. This is done using U^* to cluster the rows, during reordering, and V^* to cluster the columns, during reordering.

2.9.1 - Reorder Rows

To reorder the rows of \bar{L} , we consider \bar{L} to consist of m row vectors, denoted $\bar{l}_{v_i}^T \in \mathbb{R}^n$. That is to say that we may express the i th row of \bar{L} , $(\bar{L})_{i*}$, in terms of \bar{l}_{v_i} by the relation

$$(\bar{E})_{i*} = \vec{l}_{v_i}$$

Now, let us define $n_i \in \mathbb{Z}^+$ to be the number of elements in each of the $U_i^* \in U^*$; and $n^* \in \mathbb{Z}^+$ to be the number of U_i^* in U^* . That is to say

$$n_i = |U_i^*|$$

and

$$n^* = |U^*|$$

Using these definitions, we will reorder E , by rows, into the

temporary matrix $\bar{L} \in \mathbb{R}^{m \times n}$ via
the routine below

```

TERMLIST* = [];
t = 1;
FOR (i = 1 : u*) {
    U_i* = U*[i];
    FOR (j = 1 : u*) {
        k = U_i* [j];
        ( $\bar{L}_0$ )_{t*} = ( $\bar{L}$ )_{u*};
        TERMLIST[t] = TERMLIST[k];
        t++;
    }
}
}

```

Where the 'TERMLIST []' variable is
an array, of size m, which
contains the labels for each of the

terms corresponding to the original ordering. Similarly, the 'TERM_LIST*' variable is an array, of size m , containing the labels for each of the corresponding to the new ordering. We will next use our temporary, row-reordered matrix, \bar{E}_0 , to create the final, completely reordered matrix, E^* .

2.9.2 - Reorder Columns

Since \bar{E} and \bar{E}_0 have the same column order and \bar{E}_0 is already row-reordered, we may construct E^* from \bar{E}_0 by reordering the

columns of \bar{L}_0 according to V^* .

As with the row-reordering, we consider \bar{L}_0 to consist of n column vectors, denoted $\bar{l}_{v_j} \in \mathbb{R}^m$.

That is to say that we may express the j th column of \bar{L}_0 , $(\bar{L}_0)_{*j}$, in terms of \bar{l}_{v_j} by the relation

$$(\bar{L}_0)_{*j} = \bar{l}_{v_j}$$

Now, let us define $v_j \in \mathbb{Z}^+$ to be the number of elements in each of the $V_j \in V$; and $v^* \in \mathbb{Z}^+$ to be the number of V_j^* in V^* . That is to say that

$$v_j = |\nabla_j|$$

and

$$v^* = |\nabla^*|$$

Using these definitions, we will reorder \bar{L}_0 , by columns, into the final matrix $\bar{L}^* \in \mathbb{R}^{m \times n}$ via the routine below

```

OBJLIST* = [];
t = 1;
FOR (j = 1 : v*) {
     $\nabla_j^* = \nabla^*[j];$ 
    FOR (i = 1 : v_j) {
        k =  $\nabla_j^*[i];$  ( $\bar{L}^*$ )i,t = ( $\bar{L}_0$ )*k,t;
        OBJLIST**[t] = OBJLIST[k]; t++;
    }
}
END
}
END

```

where the 'OBJList[]' variable is an array, of size n , containing the labels for each of the objects in an order corresponding to the original column ordering of \mathbb{E} and, by extension, \mathbb{E}_0 . Similarly, the 'OBJList^t[]' variable is an array, also of size n , which contains the labels for each of the objects but ordered to correspond to the new column ordering.

4 Other Algorithms

Our algorithm will be benchmarked against two other algorithms, each with two versions. These algorithms and their versions are

- 1) SVD Signs

 - + A) of order $k < r$

 - + B) of full rank ($k=r$)

- 2) SVD Gaps

 - + A) of order $k < r$

 - + B) of full rank ($k=r$)

with $r = \text{RANK}[\mathbb{E}]$.

4.1 - SVD signs of order k

For

5 Test Data

we will use three types of data to test the algorithm.

These are

- 1) Contrived Data
- 2) Random Data
- 3) 'Real-World' Data

The contrived data will be used to ensure that the algorithm produces the expected results from a known and solved dataset. On the other hand, the random data serves to verify that this algorithm produces similar results

compared to the other algorithms
being used as benchmarks.