

Clustering by Connections in the Input Matrix

Jonathan McFadden

University of Washington Tacoma

mcfaddja@uw.edu

Summer 2016

Contents

I	Introduction, Definitions, Datasets, and Metrics	7
1	Introduction	9
2	Definitions	11
2.1	Input Matrix	11
2.2	SVD Matrices	12
2.2.1	Full SVD	12
2.2.2	Approximate SVD	12
2.3	Useful Functions	13
2.4	Connection Matrices	13
2.4.1	Row Connection Matrix	13
2.4.2	Column Connection Matrix	14
2.5	Relation Sets	15
2.5.1	Row Relation Set	15
2.5.2	Column Relation Set	16
2.6	Cluster Sets	17
2.6.1	Row Cluster Set	17
2.6.2	Column Cluster Set	17
3	Datasets	19
3.1	Contrived Data	19
3.2	Random Data	19

3.3	"Real World" Data	20
4	Metrics	21
4.1	Entropy	21
4.2	Connection Quality	21
II	Algorithms	23
5	Finding κ	25
5.1	What is κ	25
5.2	Choosing κ	27
6	SVD Signs	29
6.1	Choosing which SVD	30
6.2	Compare Row Sign Patters	30
6.3	Complete Row Cluster Sets	32
6.4	Single Sets for each Row Cluster	35
6.5	Compare Column Sign Patters	36
6.6	Complete Column Cluster Sets	38
6.7	Single Sets for each Column Cluster	39
6.8	Reordering $\bar{\mathbf{L}}$ into a new Matrix	41
6.8.1	Reorder Rows	41
6.8.2	Reorder Columns	41
6.8.3	Reorder <i>Term</i> /Row Labels	41
6.8.4	Reorder <i>Object</i> /Column Labels	41
7	SVD Gaps	43
8	CbC Disjoint Sets	45
8.1	Compute $\bar{\mathbf{A}}_U$	46
8.2	Construct \mathbb{U}	46
8.3	Complete the \mathbb{U} Sets	47

8.4	Construct the \mathbb{U}^* Sets	49
8.5	Compute $\bar{\bar{\mathbf{A}}}_V$	49
8.6	Construct \mathbb{V}	50
8.7	Complete the \mathbb{V} Sets	51
8.8	Construct the \mathbb{V}^* Sets	52
9	CbC Non-Disjoint Sets	53
III	Time Complexity Analysis	55
10	SVD Signs - Time Complexity	57
11	SVD Gaps - Time Complexity	59
12	CbC Disjoint Sets - Time Complexity	61
13	CbC Non-Disjoint Sets - Time Complexity	63
IV	Space Complexity Analysis	65
14	SVD Signs - Space Complexity	67
15	SVD Gaps - Space Complexity	69
16	CbC Disjoint Sets - Space Complexity	71
17	CbC Non-Disjoint Sets - Space Complexity	73
V	Results, Comparisons, and Conclusions	75
18	Comparison of Result Data	77
19	Comparison of Metrics on Data	79
20	Comparison of CPU Run-Time	81

21 Comparison of Memory Space Requirements	83
22 Conclusion	85

Part I

Introduction, Definitions, Datasets, and Metrics

Chapter 1

Introduction

Introduction goes here

Chapter 2

Definitions

2.1 Input Matrix

Let $\bar{\bar{\mathbf{L}}} \in \mathbb{R}^{m \times n}$ be a matrix which represents how m terms describe a collection of n objects.

We can express this formally as

$$\left(\bar{\bar{\mathbf{L}}}\right)_{ij} \equiv \begin{cases} 1 & \text{if the } i\text{th term describes the } j\text{th object.} \\ 0 & \text{otherwise} \end{cases} \quad (2.1.1)$$

We may also use the alternative definition

$$\left(\bar{\bar{\mathbf{L}}}\right)_{ij} \equiv \begin{cases} x & \text{the } i\text{th term describes the } j\text{th object } x \text{ times} \\ 0 & \text{otherwise} \end{cases} \quad (2.1.2)$$

In both definitions for the $\left(\bar{\bar{\mathbf{L}}}\right)_{ij}$, we require that $i \in [1, m] \subset \mathbb{Z}^+$ and that $j \in [1, n] \subset \mathbb{Z}^+$.

2.2 SVD Matrices

We will be using both the full SVD of $\bar{\bar{\mathbf{L}}}$ and the κ th order approximate SVD of $\bar{\bar{\mathbf{L}}}$ where the order of the approximate SVD, κ , must satisfy the relation

$$\kappa < r$$

with r defined as

$$r = \text{rank} \left[\bar{\bar{\mathbf{L}}} \right] \quad (2.2.1)$$

2.2.1 Full SVD

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

$$\bar{\bar{\mathbf{L}}} = \bar{\bar{\mathbf{U}}} \bar{\bar{\mathbf{S}}} \bar{\bar{\mathbf{V}}}^\top$$

where $\bar{\bar{\mathbf{U}}} \in \mathbb{R}^{m \times r}$, $\bar{\bar{\mathbf{S}}} \in \mathbb{R}^{r \times r}$, and $\bar{\bar{\mathbf{V}}} \in \mathbb{R}^{n \times r}$ with r as defined as above in 2.2.1. We call $\bar{\bar{\mathbf{U}}}$ the "*Row Matrix*" of $\bar{\bar{\mathbf{L}}}$, $\bar{\bar{\mathbf{S}}}$ the "*Singular Value Matrix*" of $\bar{\bar{\mathbf{L}}}$, and $\bar{\bar{\mathbf{V}}}$ the "*Column Matrix*" of $\bar{\bar{\mathbf{L}}}$.

2.2.2 Approximate SVD

Let the κ th order approximate singular value decomposition of $\bar{\bar{\mathbf{L}}}$ be expressed as

$$\bar{\bar{\mathbf{L}}} = \bar{\bar{\mathbf{U}}}_\kappa \bar{\bar{\mathbf{S}}}_\kappa \bar{\bar{\mathbf{V}}}_\kappa^\top$$

where $\bar{\bar{\mathbf{U}}}_\kappa \in \mathbb{R}^{m \times \kappa}$, $\bar{\bar{\mathbf{S}}}_\kappa \in \mathbb{R}^{\kappa \times \kappa}$, and $\bar{\bar{\mathbf{V}}}_\kappa \in \mathbb{R}^{n \times \kappa}$ with r as defined as above in 2.2.1. We call $\bar{\bar{\mathbf{U}}}_\kappa$ the "Approximate Row Matrix" or "Reduced Row Matrix" of $\bar{\bar{\mathbf{L}}}$, $\bar{\bar{\mathbf{S}}}_\kappa$ the "Approximate Singular Value Matrix" or "Reduced Singular Value Matrix" of $\bar{\bar{\mathbf{L}}}$, and $\bar{\bar{\mathbf{V}}}_\kappa$ the "Approximate Column Matrix" or "Reduced Column Matrix" of $\bar{\bar{\mathbf{L}}}$.

2.3 Useful Functions

The following function will prove useful for the alternative definitions of the *Connection Matrices* given below. This function indicates if the i th element of $\bar{\bar{\mathbf{L}}}$ can be connected to the j th element of $\bar{\bar{\mathbf{L}}}$ and is denoted by $\delta^\star(\bar{\bar{\mathbf{L}}}, i, j)$. We define this function according to the expression

$$\delta^\star(\bar{\bar{\mathbf{L}}}, i, j) \equiv \begin{cases} 1 & \text{if } \bar{\bar{\mathbf{L}}}_{ij} \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2.3.1)$$

with $i \in [1, m] \subset \mathbb{Z}^+$ and that $j \in [1, n] \subset \mathbb{Z}^+$.

2.4 Connection Matrices

We will create two "Connection Matrices" from $\bar{\bar{\mathbf{L}}}$. These matrices represent the connections among all of the row elements of $\bar{\bar{\mathbf{L}}}$ or among all of the column elements of $\bar{\bar{\mathbf{L}}}$. They will provide the information required for clustering either the rows of $\bar{\bar{\mathbf{L}}}$ or for clustering the columns of $\bar{\bar{\mathbf{L}}}$.

2.4.1 Row Connection Matrix

The first of the "Connection Matrices" that we will create from $\bar{\bar{\mathbf{L}}}$ is the "Row Connection Matrix". This matrix will be used to provide information for clustering the rows of $\bar{\bar{\mathbf{L}}}$. We

will denote this matrix by $\bar{\bar{\mathbf{A}}}_U \in \mathbb{R}^{m \times m}$ and define it, in terms of $\bar{\bar{\mathbf{L}}}$, according to the expression

$$\bar{\bar{\mathbf{A}}}_U = \bar{\bar{\mathbf{L}}} \bar{\bar{\mathbf{L}}}^\top \quad (2.4.1)$$

We may also use the alternate definition

$$\left(\bar{\bar{\mathbf{A}}}_U\right)_{ij} = \sum_{k=1}^n \left\{ \delta^\star \left(\bar{\bar{\mathbf{L}}}, i, k \right) \bar{\bar{\mathbf{L}}}_{jk} \right\} \quad (1.4.1a)$$

where $i, j \in [1, m] \subset \mathbb{Z}^+$.

2.4.2 Column Connection Matrix

The second "*Connection Matrix*" to be created from $\bar{\bar{\mathbf{L}}}$ is the "*Column Connection Matrix*". This matrix will be used to provide information for clustering the columns of $\bar{\bar{\mathbf{L}}}$.

We will denote this matrix by $\bar{\bar{\mathbf{A}}}_V \in \mathbb{R}^{n \times n}$ and define it, in terms of $\bar{\bar{\mathbf{L}}}$, according to the expression

$$\bar{\bar{\mathbf{A}}}_V = \bar{\bar{\mathbf{L}}}^\top \bar{\bar{\mathbf{L}}} \quad (2.4.2)$$

We may also use the alternate definition

$$\left(\bar{\bar{\mathbf{A}}}_V\right)_{ij} = \sum_{k=1}^m \left\{ \delta^\star \left(\bar{\bar{\mathbf{L}}}, k, i \right) \bar{\bar{\mathbf{L}}}_{kj} \right\} \quad (1.4.2a)$$

where $i, j \in [1, n] \subset \mathbb{Z}^+$.

2.5 Relation Sets

We will now define two sets of 2-tuples. The first will describe the relation or relations each row element has with all the other row elements; while the second will describe the relation or relations each column element has with all the other column elements. One of the elements of the 2-tuple represents the other rows or columns that are related to a given row or column. The other element of the 2-tuple describes the strength of those relations.

2.5.1 Row Relation Set

Let \mathbb{U} be a set of m 2-tuples that can be expressed as

$$\mathbb{U} = \left\{ \{\mathbb{U}_1, \mathbb{M}_1\}, \{\mathbb{U}_2, \mathbb{M}_2\}, \dots, \{\mathbb{U}_m, \mathbb{M}_m\} \right\}$$

where the $\{\mathbb{U}_i, \mathbb{M}_i\}$ are 2-tuples where the first element of the tuple, \mathbb{U}_i , represents which rows in $\bar{\bar{\mathbf{L}}}$ are connected to the i th row of $\bar{\bar{\mathbf{L}}}$; while the second element of the tuple, \mathbb{M}_i , describes the strength of those relations, with i such that $i \in [1, m] \subset \mathbb{Z}^+$. We now define \mathbb{U}_i and \mathbb{M}_i concurrently as

$$\mathbb{U}_i \equiv \left\{ j \mid j \in [1, m] \subset \mathbb{Z}^+ \text{ and } \left(\bar{\bar{\mathbf{A}}}_U \right)_{ij} \neq 0 \right\} \quad (2.5.1)$$

and

$$\mathbb{M}_i \equiv \left\{ \left(\bar{\bar{\mathbf{A}}}_U \right)_{ij} \mid j \in [1, m] \subset \mathbb{Z}^+ \text{ and } \left(\bar{\bar{\mathbf{A}}}_U \right)_{ij} \neq 0 \right\} \quad (2.5.2)$$

respectively. These definitions require that the expression

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

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

2.5.2 Column Relation Set

Let \mathbb{V} be a set of n 2-tuples that can be expressed as

$$\mathbb{V} = \left\{ \{\mathbb{V}_1, \mathbb{N}_1\}, \{\mathbb{V}_2, \mathbb{N}_2\}, \dots, \{\mathbb{V}_n, \mathbb{N}_n\} \right\}$$

where the $\{\mathbb{V}_j, \mathbb{N}_j\}$ are 2-tuples where the first element of the tuple, \mathbb{V}_j , represents which columns in $\bar{\mathbf{L}}$ are connected to the j th column of $\bar{\mathbf{L}}$; while the second element of the tuple, \mathbb{N}_j , describes the strength of those relations, with j such that $j \in [1, n] \subset \mathbb{Z}^+$. We now define \mathbb{V}_j and \mathbb{N}_j concurrently as

$$\mathbb{V}_j \equiv \left\{ i \mid i \in [1, n] \subset \mathbb{Z}^+ \text{ and } \left(\bar{\mathbf{A}}_V \right)_{ji} \neq 0 \right\} \quad (2.5.3)$$

and

$$\mathbb{N}_j \equiv \left\{ \left(\bar{\mathbf{A}}_V \right)_{ji} \mid i \in [1, n] \subset \mathbb{Z}^+ \text{ and } \left(\bar{\mathbf{A}}_V \right)_{ji} \neq 0 \right\} \quad (2.5.4)$$

respectively. These definitions require that the expression

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

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

2.6 Cluster Sets

We now define two sets of sets with the first set of sets representing the row clusters and the second set of sets representing the column clusters. Each set of sets is composed of sets which represent the members of each row or column cluster.

2.6.1 Row Cluster Set

The "*Row Cluster Set*" is a set of $\mathbf{u} \in \mathbb{Z}^+$ sets with each component set representing one of the \mathbf{u} row clusters. We will denote this set of sets as \mathbb{U}^* and formally describe it as

$$\mathbb{U}^* = \left\{ \mathbb{U}_1^*, \mathbb{U}_2^*, \dots, \mathbb{U}_{\mathbf{u}}^* \right\}$$

where the \mathbb{U}_μ^* are sets which contain the members of the μ th row cluster, with $\mu \in [1, \mathbf{u}] \subset \mathbb{Z}^+$. Each of the \mathbb{U}_μ^* is constructed from the $\{\mathbb{U}_i, \mathbb{M}_i\}$ of its (the \mathbb{U}_μ^*) constituent elements via a process we will describe later.

2.6.2 Column Cluster Set

The "*Column Cluster Set*" is a set of $\mathbf{v} \in \mathbb{Z}^+$ sets with each component set representing one of the \mathbf{v} column clusters. We will denote this set of sets as \mathbb{V}^* and formally describe it as

$$\mathbb{V}^\star = \left\{ \mathbb{V}_1^\star, \mathbb{V}_2^\star, \dots, \mathbb{V}_{\mathfrak{v}}^\star \right\}$$

where the \mathbb{V}_ν^\star are sets which contain the members of the ν th column cluster, with $\nu \in [1, \mathfrak{v}] \subset \mathbb{Z}^+$. Each of the \mathbb{V}_ν^\star is constructed from the $\{\mathbb{V}_j, \mathbb{N}_j\}$ of its (the \mathbb{V}_ν^\star) constituent elements via a process we will describe later.

Chapter 3

Datasets

Initial Dataset info goes here

There 3 types of datasets

1. Contrived Datasets
2. Random Datasets
3. "*Real World*" Datasets

The "*Real World*" will initially be truncated due to limits on processing and memory; however, later, a few of the "*Real World*" datasets will be run in full.

3.1 Contrived Data

Contrived Dataset info goes here

3.2 Random Data

Random Dataset info goes here

3.3 ”*Real World*” Data

”*Real World*” Dataset info goes here

Chapter 4

Metrics

Initial Metrics info goes here

4.1 Entropy

Initial Entropy Metric info goes here

4.2 Connection Quality

Initial Connection Quality Metric info goes here

Part II

Algorithms

Chapter 5

Finding κ

Both the SVD Signs Algorithm and the SVD Gaps Algorithm require a value for κ when we use the Reduced SVD (or Approximate SVD) of $\bar{\bar{\mathbf{L}}}$. Thus, we must have both an understanding of what κ represents, as well as a process for reliably determining its value.

5.1 What is κ

To begin to understand what κ is, let us define $r \in \mathbb{Z}^+$ as

$$r = \text{rank} \left[\bar{\bar{\mathbf{L}}} \right]$$

Next, recall the expression for the full Singular Value Decomposition of $\bar{\bar{\mathbf{L}}}$ from chapter 2 written below

$$\bar{\bar{\mathbf{L}}} = \bar{\bar{\mathbf{U}}} \bar{\bar{\mathbf{S}}} \bar{\bar{\mathbf{V}}}^\top,$$

where $\bar{\bar{\mathbf{U}}} \in \mathbb{R}^{m \times r}$, $\bar{\bar{\mathbf{S}}} \in \mathbb{R}^{r \times r}$, and $\bar{\bar{\mathbf{V}}} \in \mathbb{R}^{n \times r}$. This is expensive to compute, in both time and

memory; therefore, we seek something which can be computed for less cost.

Let us now create new versions of the $\bar{\mathbf{U}}$, $\bar{\mathbf{S}}$, and $\bar{\mathbf{V}}$ matrices, denoted $\bar{\bar{\mathbf{U}}}_\xi$, $\bar{\bar{\mathbf{S}}}_\xi$, and $\bar{\bar{\mathbf{V}}}_\xi$, respectively. We will define these new versions of the $\bar{\mathbf{U}}$, $\bar{\mathbf{S}}$, and $\bar{\mathbf{V}}$ matrices as such that

$$\bar{\bar{\mathbf{U}}}_\xi \in \mathbb{R}^{m \times \xi}$$

$$\bar{\bar{\mathbf{S}}}_\xi \in \mathbb{R}^{\xi \times \xi}$$

and

$$\bar{\bar{\mathbf{V}}}_\xi \in \mathbb{R}^{n \times \xi}$$

where $\xi \in \mathbb{Z}^{++} \ni \xi < r$. Since $\xi < r$, these new matrices are obviously less expensive to store. Additionally, note that if we compute the matrix product of $\bar{\bar{\mathbf{U}}}_\xi$, $\bar{\bar{\mathbf{S}}}_\xi$, and $\bar{\bar{\mathbf{V}}}_\xi$, we obtain the result

$$\left(\bar{\bar{\mathcal{L}}} = \bar{\bar{\mathbf{U}}}_\xi \bar{\bar{\mathbf{S}}}_\xi \bar{\bar{\mathbf{V}}}_\xi^\top \right) \in \mathbb{R}^{m \times n} \quad (5.1.1)$$

which is clearly such that $\bar{\bar{\mathcal{L}}} \in \mathbb{R}^{m \times n}$. Since $\bar{\bar{\mathbf{S}}}$ is the *Singular Value Matrix* of $\bar{\bar{\mathbf{L}}}$, its only non-zero elements are the r singular values of $\bar{\bar{\mathbf{L}}}$. These elements are arranged along the diagonal of $\bar{\bar{\mathbf{S}}}$ and ordered from highest to lowest. That is to say, for the singular values of $\bar{\bar{\mathbf{L}}}$, $\{\sigma_1, \sigma_2, \dots, \sigma_r\}$ we have $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_i \geq \dots \geq \sigma_r$. Now, the singular values of $\bar{\bar{\mathbf{L}}}$ can be found by computing the eigenvalues of $\bar{\bar{\mathbf{L}}}^\top \bar{\bar{\mathbf{L}}}$ (denoted λ_i^*), followed by taking the real component of the square root of each of these eigenvalues. That is to say, the value of each σ_i can be found using the relation

$$\sigma_i = \text{Re} \left[\sqrt{\lambda_i^*} \right] \quad (5.1.2)$$

Thus the diagonal elements of $\bar{\bar{\mathbf{S}}}$, $\left\{ \left(\bar{\bar{\mathbf{S}}} \right)_{11}, \left(\bar{\bar{\mathbf{S}}} \right)_{22}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{ii}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{rr} \right\}$, are such that

$$\left(\bar{\bar{\mathbf{S}}} \right)_{11} \geq \left(\bar{\bar{\mathbf{S}}} \right)_{22} \geq \dots \geq \left(\bar{\bar{\mathbf{S}}} \right)_{ii} \geq \dots \geq \left(\bar{\bar{\mathbf{S}}} \right)_{rr}$$

Since our new $\bar{\bar{\mathbf{S}}}_\xi$ is constructed from $\bar{\bar{\mathbf{S}}}$, $\bar{\bar{\mathbf{S}}}_\xi$ also contains the singular values of $\bar{\bar{\mathbf{L}}}$ arranged along its diagonal by magnitude. Therefore, the elements of $\bar{\bar{\mathbf{S}}}_\xi$ must be such that

$$\left\{ \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{11}, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{22}, \dots, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{\iota\iota}, \dots, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{\xi\xi} \right\} \subset \left\{ \left(\bar{\bar{\mathbf{S}}} \right)_{11}, \left(\bar{\bar{\mathbf{S}}} \right)_{22}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{ii}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{rr} \right\}$$

Additionally, since $\xi < r$, we have

$$\left| \left\{ \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{11}, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{22}, \dots, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{\iota\iota}, \dots, \left(\bar{\bar{\mathbf{S}}}_\xi \right)_{\xi\xi} \right\} \right| < \left| \left\{ \left(\bar{\bar{\mathbf{S}}} \right)_{11}, \left(\bar{\bar{\mathbf{S}}} \right)_{22}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{ii}, \dots, \left(\bar{\bar{\mathbf{S}}} \right)_{rr} \right\} \right|$$

All that remains is to decide *which* elements of $\bar{\bar{\mathbf{S}}}$ to include in $\bar{\bar{\mathbf{S}}}_\xi$, subject to the criteria above. Let us

5.2 Choosing κ

From here, there are six possible ways to choose κ . The first five are quite simple and are

1. $\kappa \ni \forall i \in [1, k], \sigma_i > (\sigma_1/2)$
2. $\kappa \ni \forall i \in [1, k], \sigma_i > \sqrt{\sigma_1}$

3. $\kappa \ni \forall i \in [1, k], \sigma_i > \left(\frac{1}{2}((\sigma_1/2) + \sqrt{\sigma_1})\right)$
4. $\kappa \ni \forall i \in [1, k], \sigma_i > \min[(\sigma_1/2), \sqrt{\sigma_1}]$
5. $\kappa \ni \forall i \in [1, k], \sigma_i > \max[(\sigma_1/2), \sqrt{\sigma_1}]$

The sixth option is considerably more complicated as it requires find the point (singular value index) where the rate of change in the difference between adjacent singular values, $\Delta\sigma_i$, changes most sharply. Where we define $\Delta\sigma_i$ as

$$\Delta\sigma_i = \begin{cases} 0 & \text{for } i = 1 \\ \sigma_i - \sigma_{i+1} & \text{otherwise} \end{cases} \quad (5.2.1)$$

This is to say, we seek the point at which the difference between adjacent $\Delta\sigma_i$ changes most abruptly. Thus, we seek the value for j , where $j \in (1, r]$, such that

$$\Delta\sigma_j - \Delta\sigma_{j+1}$$

is a maximum. Graphically, this appears as an 'elbow' in the graph of the singular values versus their associated indices. The index, j , at which this 'elbow' occurs represents the desired value for κ . If time permits, we will attempt to create an algorithm to automatically find the value for κ using this method; however, for now, this will have to be done manually.

Finally, with the required value of κ in hand, it is possible to proceed to each of the three clustering algorithm that will be employed, starting with the SVD Signs algorithm.

Chapter 6

SVD Signs

The SVD Signs Algorithm has algorithm has eight steps of its own and one step it shares with the SVD Gaps Algorithm. The step shared between the These steps are

1. Compute the k th order approximate SVD of $\bar{\bar{\mathbf{A}}}$.
2. Compare the row sign patterns.
3. Complete the row cluster's sets.
4. Create single sets for each row cluster.
5. Compare the column sign patterns.
6. Complete the column cluster's sets.
7. Create single sets for each column cluster.
8. Reorder $\bar{\bar{\mathbf{A}}}$ so that rows and columns from the same cluster have indices which are adjacent.

6.1 Choosing which SVD

For the first step, we simply use the expression from (1.2) with the value for k found previously. Thus, we will proceed with

$$\bar{\bar{\mathbf{A}}} \approx \bar{\bar{\mathbf{U}}}_k \bar{\bar{\mathbf{S}}}_k \bar{\bar{\mathbf{V}}}_k^T \quad (1.2)$$

from which we will use $\bar{\bar{\mathbf{U}}}_k$ for clustering rows and $\bar{\bar{\mathbf{V}}}_k$ for clustering columns.

Choosing which SVD section SVD Signs Algorithm info goes here

6.2 Compare Row Sign Patterns

To simplify the comparison of sign patterns between the rows of $\bar{\bar{\mathbf{U}}}_k$, we will create a new matrix $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})} \in \mathbb{R}^{m \times k}$. This new matrix will be based on $\bar{\bar{\mathbf{U}}}_k$, with the values for the elements of $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}$ being determined by the following definition

$$\left(\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}\right)_{ij} = \begin{cases} 1 & \text{if } \left(\bar{\bar{\mathbf{U}}}_k\right)_{ij} > 0 \\ 0 & \text{if } \left(\bar{\bar{\mathbf{U}}}_k\right)_{ij} = 0 \\ -1 & \text{if } \left(\bar{\bar{\mathbf{U}}}_k\right)_{ij} < 0 \end{cases} \quad (6.2.1)$$

Here, the $\left(\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}\right)_{ij}$ are the elements of $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}$, the $\left(\bar{\bar{\mathbf{U}}}_k\right)_{ij}$ are the elements of $\bar{\bar{\mathbf{U}}}_k$, and, for both $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}$ and $\left(\bar{\bar{\mathbf{U}}}_k\right)_{ij}$, the indices i and j are such that $i \in [1, m]$ and $j \in [1, k]$.

Next, we will create a new m by m matrix to represent the connections between the rows of $\bar{\bar{\mathbf{U}}}_k$. This new matrix will be denoted by $\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}$, with $\bar{\bar{\mathbf{U}}}_k^{(\text{conn})} \in \mathbb{R}^{m \times m}$ and having

elements $\left(\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}\right)_{ij}$, for $i, j \in [1, m]$. We will define each $\left(\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}\right)_{ij}$ to be 1 if the i th row of $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}$ is equivalent to the j th row of $\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}$. That is to say, formally, that

$$\left(\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}\right)_{ij} = \begin{cases} 1 & \text{if } \forall l \in [1, k], \left(\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}\right)_{il} = \left(\bar{\bar{\mathbf{U}}}_k^{(\text{sign})}\right)_{jl} \text{ holds} \\ 0 & \text{otherwise} \end{cases} \quad (6.2.2)$$

Additionally, we note that the definition in (1.7) implies that the relation

$$\left(\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}\right)_{ii} = 1, \forall i \in [1, m]$$

must be valid as well.

The final step in comparing the sign patterns of the rows in $\bar{\bar{\mathbf{U}}}_k$ is to collect the connections between rows. These connections between rows, indicated by their sign patterns, will be collected into the set of sets $\mathbb{U}_k^{(\text{rel})}$ with m elements such that $\mathbb{U}_k^{(\text{rel})}$ may be defined

$$\mathbb{U}_k^{(\text{rel})} = \left\{ \left(\mathbb{U}_k^{(\text{rel})}\right)_1, \left(\mathbb{U}_k^{(\text{rel})}\right)_2, \dots, \left(\mathbb{U}_k^{(\text{rel})}\right)_m \right\}$$

where the elements of $\mathbb{U}_k^{(\text{rel})}$ are all sets in their own right and are denoted by the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ for $i \in [1, m]$. Each of these k elements in $\mathbb{U}_k^{(\text{rel})}$ is defined such that for $\forall i \in [1, m]$ the i th element, $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$, is the set

$$\left(\mathbb{U}_k^{(\text{rel})}\right)_i \equiv \left\{ j \mid j \in [1, m] \in \left(\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}\right)_{ij} \neq 0 \right\}, \forall i \in [1, k] \quad (6.2.3)$$

Additionally, note that this definition must imply that $\{i\} \in \left(\mathbb{U}_k^{(\text{rel})}\right)_i$, $\forall i \in [1, m]$ is valid as well. The elements of each of the sets $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ as defined above should only be considered as the initial elements of each $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$. This is due to the fact that additional elements may be added to each $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ as described in the next section.

6.3 Complete Row Cluster Sets

The elements of $\mathbb{U}_k^{(\text{rel})}$ should contain the indices of all rows in the same cluster as each element. That is to say, each $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ should be the set containing the indices of all the rows in the same cluster as row i , in addition to the index of row i ; however, this is not guaranteed to be the case initially after the creation of the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ as described above. To illustrate this, consider the values for some $\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}$ as given below

$$\bar{\bar{\mathbf{U}}}_k^{(\text{conn})} = \begin{bmatrix} 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix} \quad (6.3.1)$$

Using the definition for the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ from (1.8), we obtain the following initial values for the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ generated from the $\bar{\bar{\mathbf{U}}}_k^{(\text{conn})}$ expressed above

i	$\mathbb{U}_k^{(\text{rel})}$
1	$\left(\mathbb{U}_k^{(\text{rel})}\right)_1 = \{1, 2\}$
2	$\left(\mathbb{U}_k^{(\text{rel})}\right)_2 = \{2\}$
3	$\left(\mathbb{U}_k^{(\text{rel})}\right)_3 = \{2, 3\}$
4	$\left(\mathbb{U}_k^{(\text{rel})}\right)_4 = \{4\}$
5	$\left(\mathbb{U}_k^{(\text{rel})}\right)_5 = \{5\}$
6	$\left(\mathbb{U}_k^{(\text{rel})}\right)_6 = \{5, 6\}$

Table ??: Example of initial values of $\mathbb{U}_k^{(\text{rel})}$

Clearly, rows 1, 2, and 3 belong in the same cluster; rows 5 and 6 also belong to the same cluster, but one different than the first; and row 4 belongs by itself. However, with the exception of rows 4 and 6, the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ do not reflect this, at least not initially. Thus, each of the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ must be completed by adding the other indices of rows which are in the same cluster as the row with which $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ is associated. This will be accomplished using the intersect, $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{U}_k^{(\text{rel})}\right)_j$, and union, $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{U}_k^{(\text{rel})}\right)_j$, set operations.

We will use the intersect set operation will indicate if the rows represented by $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ and $\left(\mathbb{U}_k^{(\text{rel})}\right)_j$ belong to the same cluster, as the operation will give

$$\left(\mathbb{U}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{U}_k^{(\text{rel})}\right)_j = \emptyset$$

if rows i and j belong to different clusters. If this intersection yields a non-empty set, then we will use the second set operation, the union

$$\left(\mathbb{U}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{U}_k^{(\text{rel})}\right)_j$$

to combine the elements of $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ and $\left(\mathbb{U}_k^{(\text{rel})}\right)_j$. This discussion allows us to layout the following process for completing all of the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \in \mathbb{U}_k^{(\text{rel})}$.

Data: Initial values for the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \in \mathbb{U}_k^{(\text{rel})}$

Result: Completed values for the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \in \mathbb{U}_k^{(\text{rel})}$, fully representing the cluster of each row by containing all elements of that cluster.

```

begin
  for  $i = 1 : m$  do
    for  $j = 1 : m$  do
      if  $\left(\mathbb{U}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{U}_k^{(\text{rel})}\right)_j \neq \emptyset$  then
         $\left(\mathbb{U}_k^{(\text{rel})}\right)_i = \left(\mathbb{U}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{U}_k^{(\text{rel})}\right)_j$  ;
         $\left(\mathbb{U}_k^{(\text{rel})}\right)_j = \left(\mathbb{U}_k^{(\text{rel})}\right)_j \cup \left(\mathbb{U}_k^{(\text{rel})}\right)_i$  ;
      end
    end
  end
end
end

```

Algorithm 1: Completing the sets $\left(\mathbb{U}_k^{(\text{rel})}\right)_i, \forall i \in [i, m]$

Applying this algorithm to the example in **Table ??**, we have

i	$\mathbb{U}_k^{(\text{rel})}$	$i = 1$ $j = 2$	$i = 1$ $j = 3$	$i = 2$ $j = 2$	$i = 5$ $j = 6$
1	$\left(\mathbb{U}_k^{(\text{rel})}\right)_1 = \{1, 2\}$		$= \{1, 2, 3\}$		
2	$\left(\mathbb{U}_k^{(\text{rel})}\right)_2 = \{2\}$	$= \{1, 2\}$		$= \{1, 2, 3\}$	
3	$\left(\mathbb{U}_k^{(\text{rel})}\right)_3 = \{2, 3\}$		$= \{1, 2, 3\}$		
4	$\left(\mathbb{U}_k^{(\text{rel})}\right)_4 = \{4\}$				
5	$\left(\mathbb{U}_k^{(\text{rel})}\right)_5 = \{5\}$				$= \{5, 6\}$
6	$\left(\mathbb{U}_k^{(\text{rel})}\right)_6 = \{5, 6\}$				

Table ??: Example of algorithm on the elements of $\mathbb{U}_k^{(rel)}$,
with steps and elements without change omitted.

6.4 Single Sets for each Row Cluster

Now that each set in $\mathbb{U}_k^{(rel)}$ contains all of the elements in the same cluster as the row represented by that set, the next step is to uniquely identify each cluster by numbering them. Additionally, we wish to associate the index number for each cluster with the set listing all of the elements of the cluster in question. To do this, we create a new set of sets, $\mathbb{U}_k^{(clust)}$, defined as

$$\mathbb{U}_k^{(clust)} = \left\{ \left(\mathbb{U}_k^{(clust)} \right)_1, \left(\mathbb{U}_k^{(clust)} \right)_2, \dots, \left(\mathbb{U}_k^{(clust)} \right)_c \right\}$$

whose elements are sets of 2-tuples and with its cardinality, $|\mathbb{U}_k^{(clust)}| = c$, being the number of unique clusters. These sets of 2-tuples that comprise $\mathbb{U}_k^{(clust)}$ are denoted $\left(\mathbb{U}_k^{(clust)} \right)_i$ and are defined

$$\left(\mathbb{U}_k^{(clust)} \right)_i = \left\{ i, \left(\mathbb{U}_k^{(rel)} \right)_j \right\} \quad (6.4.1)$$

The first element of these 2-tuples, i , is the index of the cluster in question and the second element of these 2-tuples, $\left(\mathbb{U}_k^{(rel)} \right)_j$, is the set of all elements in that cluster, with the index j being the lowest index of $\mathbb{U}_k^{(rel)}$ where that set occurs. To build $\mathbb{U}_k^{(clust)}$, we use the following process

Later, we will use the elements of $\mathbb{U}_k^{(clust)}$ generated by this process to reorder the rows of $\bar{\mathbf{A}}$. However, we will next repeat the proceeding three steps for each of the n columns in $\bar{\mathbf{A}}$.

Data: Initial values for the $(\mathbb{U}_k^{(\text{rel})})_i \in \mathbb{U}_k^{(\text{rel})}$

Result: Completed $\mathbb{U}_k^{(\text{clust})}$.

```

begin
   $(\mathbb{U}_k^{(\text{clust})})_1 \leftarrow \{1, (\mathbb{U}_k^{(\text{rel})})_1\};$ 
   $c_0 \leftarrow 1;$ 
  for  $i = 2 : m$  do
    for  $j = 1 : c_0$  do
      if  $(\mathbb{U}_k^{(\text{rel})})_i \cap (\mathbb{U}_k^{(\text{clust})})_j = \emptyset$  then
         $c_0 \leftarrow (c_0 + 1);$ 
         $(\mathbb{U}_k^{(\text{clust})})_{c_0} \leftarrow \{c_0, (\mathbb{U}_k^{(\text{rel})})_i\};$ 
      end
    end
  end
end

```

Algorithm 2: Creating the sets $(\mathbb{U}_k^{(\text{clust})})_i \in \mathbb{U}_k^{(\text{clust})}, \forall i \in [i, c]$.

6.5 Compare Column Sign Patterns

To simplify the comparison of sign patterns between the columns of $\bar{\bar{\mathbf{V}}}_k^\top = \bar{\bar{\mathcal{V}}}_k$, we will create a new matrix $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})} \in \mathbb{R}^{k \times n}$. This new matrix will be based on $\bar{\bar{\mathcal{V}}}_k$, with the values for the elements of $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}$ being determined by the following definition

$$(\bar{\bar{\mathcal{V}}}_k^{(\text{sign})})_{ij} = \begin{cases} 1 & \text{if } (\bar{\bar{\mathbf{V}}}_k^\top)_{ij} > 0 \\ 0 & \text{if } (\bar{\bar{\mathbf{V}}}_k^\top)_{ij} = 0 \\ -1 & \text{if } (\bar{\bar{\mathbf{V}}}_k^\top)_{ij} < 0 \end{cases} \quad (6.5.1)$$

Here, the $(\bar{\bar{\mathcal{V}}}_k^{(\text{sign})})_{ij}$ are the elements of $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}$, the $(\bar{\bar{\mathbf{V}}}_k^\top)_{ij}$ are the elements of $\bar{\bar{\mathbf{V}}}_k^\top$, and, for both $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}$ and $(\bar{\bar{\mathbf{V}}}_k^\top)_{ij}$, the indices i and j are such that $i \in [1, k]$ and $j \in [1, n]$.

Next, we will create a new n by n matrix to represent the connections between the rows of $\bar{\bar{\mathbf{V}}}_k^\top$. This new matrix will be denoted by $\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}$, with $\bar{\bar{\mathcal{V}}}_k^{(\text{conn})} \in \mathbb{R}^{n \times n}$ and having elements

$\left(\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}\right)_{ij}$, for $i, j \in [1, n]$. We will define each $\left(\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}\right)_{ij}$ to be 1 if the i th row of $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}$ is equivalent to the j th row of $\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}$. That is to say, formally, that

$$\left(\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}\right)_{ij} = \begin{cases} 1 & \text{if } \forall l \in [1, k], \left(\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}\right)_{li} = \left(\bar{\bar{\mathcal{V}}}_k^{(\text{sign})}\right)_{lj} \quad \underline{\text{holds}} \\ 0 & \underline{\text{otherwise}} \end{cases} \quad (6.5.2)$$

Additionally, we note that the definition in (1.7) implies that the relation

$$\left(\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}\right)_{ii} = 1, \forall i \in [1, n]$$

must be valid as well.

The final step in comparing the sign patterns of the rows in $\bar{\bar{\mathbf{V}}}_k^T$ is to collect the connections between rows. These connections between rows, indicated by their sign patterns, will be collected into the set of sets $\mathbb{V}_k^{(\text{rel})}$ with m elements such that $\mathbb{V}_k^{(\text{rel})}$ may be defined

$$\mathbb{V}_k^{(\text{rel})} = \left\{ \left(\mathbb{V}_k^{(\text{rel})}\right)_1, \left(\mathbb{V}_k^{(\text{rel})}\right)_2, \dots, \left(\mathbb{V}_k^{(\text{rel})}\right)_n \right\}$$

where the elements of $\mathbb{V}_k^{(\text{rel})}$ are all sets in their own right and are denoted by the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ for $i \in [1, n]$. Each of these k elements in $\mathbb{V}_k^{(\text{rel})}$ is defined such that for $\forall i \in [1, n]$ the i th element, $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$, is the set

$$\left(\mathbb{V}_k^{(\text{rel})}\right)_i \equiv \left\{ j \mid j \in [1, m] \in \left(\bar{\bar{\mathcal{V}}}_k^{(\text{conn})}\right)_{ij} \neq 0 \right\}, \forall i \in [1, k] \quad (6.5.3)$$

Additionally, note that this definition must imply that $\{i\} \in \left(\mathbb{V}_k^{(\text{rel})}\right)_i$, $\forall i \in [1, n]$ is valid as well. The elements of each of the sets $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ as defined above should only be considered as the initial elements of each $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$. This is due to the fact that additional elements may be added to each $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ as described in the next section.

6.6 Complete Column Cluster Sets

The elements of $\mathbb{V}_k^{(\text{rel})}$ should contain the indices of all columns in the same cluster as each element. That is to say, each $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ should be the set containing the indices of all the columns in the same cluster as column i , in addition to the index of column i ; however, this is not guaranteed to be the case initially after the creation of the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ as described above. The same illustration we gave for the $\left(\mathbb{U}_k^{(\text{rel})}\right)_i$ **Section 1.2.2** illustrates this point as well. Similar to the case with the rows, each of the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ must be completed by adding the other indices of columns which are in the same cluster as the column with which $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ is associated. This will be accomplished using the intersection, $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{V}_k^{(\text{rel})}\right)_j$, and union, $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{V}_k^{(\text{rel})}\right)_j$, set operations.

We will use the intersect set operation will indicate if the columns represented by $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ and $\left(\mathbb{V}_k^{(\text{rel})}\right)_j$ belong to the same cluster, as the operation will give

$$\left(\mathbb{V}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{V}_k^{(\text{rel})}\right)_j = \emptyset$$

if columns i and j belong to different clusters. If this intersection yields a non-empty set, then we will use the second set operation, the union

$$\left(\mathbb{V}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{V}_k^{(\text{rel})}\right)_j$$

to combine the elements of $\left(\mathbb{V}_k^{(\text{rel})}\right)_i$ and $\left(\mathbb{V}_k^{(\text{rel})}\right)_j$. This discussion allows us to layout the following process for completing all of the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \in \mathbb{V}_k^{(\text{rel})}$.

Data: Initial values for the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \in \mathbb{V}_k^{(\text{rel})}$

Result: Completed values for the $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \in \mathbb{V}_k^{(\text{rel})}$, fully representing the cluster of each column by containing all elements of that cluster.

```

begin
  for i = 1 : n do
    for j = 1 : n do
      if  $\left(\mathbb{V}_k^{(\text{rel})}\right)_i \cap \left(\mathbb{V}_k^{(\text{rel})}\right)_j \neq \emptyset$  then
         $\left(\mathbb{V}_k^{(\text{rel})}\right)_i = \left(\mathbb{V}_k^{(\text{rel})}\right)_i \cup \left(\mathbb{V}_k^{(\text{rel})}\right)_j$  ;
         $\left(\mathbb{V}_k^{(\text{rel})}\right)_j = \left(\mathbb{V}_k^{(\text{rel})}\right)_j \cup \left(\mathbb{V}_k^{(\text{rel})}\right)_i$  ;
      end
    end
  end
end

```

Algorithm 3: Completing the sets $\left(\mathbb{V}_k^{(\text{rel})}\right)_i, \forall i \in [i, n]$

6.7 Single Sets for each Column Cluster

Now that each set in $\mathbb{V}_k^{(\text{rel})}$ contains all of the elements in the same cluster as the column represented by that set, the next step is to uniquely identify each cluster by numbering them. Additionally, we wish to associate the index number for each cluster with the set listing all of the elements of the cluster in question. To do this, we create a new set of sets, $\mathbb{V}_k^{(\text{clust})}$, defined as

$$\mathbb{V}_k^{(\text{clust})} = \left\{ \left(\mathbb{V}_k^{(\text{clust})} \right)_1, \left(\mathbb{V}_k^{(\text{clust})} \right)_2, \dots, \left(\mathbb{V}_k^{(\text{clust})} \right)_c \right\}$$

whose elements are sets of 2-tuples and with its cardinality, $\left| \mathbb{V}_k^{(\text{clust})} \right| = c$, being the number of unique clusters. These sets of 2-tuples that comprise $\mathbb{V}_k^{(\text{clust})}$ are denoted $\left(\mathbb{V}_k^{(\text{clust})} \right)_i$ and are defined

$$\left(\mathbb{V}_k^{(\text{clust})} \right)_i = \left\{ i, \left(\mathbb{V}_k^{(\text{rel})} \right)_j \right\} \quad (6.7.1)$$

The first element of these 2-tuples, i , is the index of the cluster in question and the second element of these 2-tuples, $\left(\mathbb{V}_k^{(\text{rel})} \right)_j$, is the set of all elements in that cluster, with the index j being the lowest index of $\mathbb{V}_k^{(\text{rel})}$ where that set occurs. To build $\mathbb{V}_k^{(\text{clust})}$, we use the following process

Data: Initial values for the $\left(\mathbb{V}_k^{(\text{rel})} \right)_i \in \mathbb{V}_k^{(\text{rel})}$
Result: Completed $\mathbb{V}_k^{(\text{clust})}$.
begin
 $\left(\mathbb{V}_k^{(\text{clust})} \right)_1 \leftarrow \left\{ 1, \left(\mathbb{V}_k^{(\text{rel})} \right)_1 \right\};$
 $c_0 \leftarrow 1;$
 for $i = 2 : n$ **do**
 for $j = 1 : c_0$ **do**
 if $\left(\mathbb{V}_k^{(\text{rel})} \right)_i \cap \left(\mathbb{V}_k^{(\text{clust})} \right)_j = \emptyset$ **then**
 $c_0 \leftarrow (c_0 + 1);$
 $\left(\mathbb{V}_k^{(\text{clust})} \right)_{c_0} \leftarrow \left\{ c_0, \left(\mathbb{V}_k^{(\text{rel})} \right)_i \right\};$
 end
 end
 end
end

Algorithm 4: Creating the sets $\left(\mathbb{V}_k^{(\text{clust})} \right)_i \in \mathbb{V}_k^{(\text{clust})}, \forall i \in [i, c]$.

6.8 Reordering $\bar{\bar{\mathbf{L}}}$ into a new Matrix

The last step is to finally reorder $\bar{\bar{\mathbf{L}}}$ into the clustered matrix $\bar{\bar{\mathbf{L}}}^* \in \mathbb{R}^{m \times n}$, which will be our result. To do this, we first use the information about the row clusters contained in $\mathbb{U}^{(\text{clust})}$ or $\mathbb{U}_k^{(\text{clust})}$ to reorder the rows of $\bar{\bar{\mathbf{L}}}$ into $\bar{\bar{\mathbf{L}}}^*$. Similarly, we use the information about the column clusters contained in $\mathbb{V}^{(\text{clust})}$ or $\mathbb{V}_k^{(\text{clust})}$ to reorder the columns of $\bar{\bar{\mathbf{L}}}$ into $\bar{\bar{\mathbf{L}}}^*$.

Additionally, we must use the information in $\mathbb{U}^{(\text{clust})}$ or $\mathbb{U}_k^{(\text{clust})}$ to ensure that the *term* label associated with each row (*representing Terms*) in $\bar{\bar{\mathbf{L}}}$ is moved so that its position corresponds to the new position of its associated row in $\bar{\bar{\mathbf{L}}}^*$. In similar fashion, the information in $\mathbb{V}^{(\text{clust})}$ or $\mathbb{V}_k^{(\text{clust})}$ is used to ensure that the *object* label associated with each column (*representing Objects*) in $\bar{\bar{\mathbf{L}}}$ is moved so that its position corresponds to the new position of its associated column in $\bar{\bar{\mathbf{L}}}^*$.

6.8.1 Reorder Rows

Process for reordering rows goes here.

6.8.2 Reorder Columns

Process for reordering columns goes here.

6.8.3 Reorder *Term*/Row Labels

Process for reordering *Term*/Row Labels goes here.

6.8.4 Reorder *Object*/Column Labels

Process for reordering *Object*/Column Labels goes here.

Chapter 7

SVD Gaps

Intro SVD Gaps Algorithm info goes here

Chapter 8

CbC Disjoint Sets

When $\bar{\bar{\mathbf{L}}}$ consists of a series of disjoint sets, our clustering algorithm proceeds according to the following simple steps

1. Compute $\bar{\bar{\mathbf{A}}}_U$
2. Construct $\mathbb{U} = \left\{ \{\mathbb{U}_1, \mathbb{M}_1\}, \{\mathbb{U}_2, \mathbb{M}_2\}, \dots, \{\mathbb{U}_m, \mathbb{M}_m\} \right\}$
3. Complete the $\{\mathbb{U}_i, \mathbb{M}_i\} \in \mathbb{U}$
4. Construct the \mathbb{U}^* set of sets
5. Compute $\bar{\bar{\mathbf{A}}}_V$
6. Construct $\mathbb{V} = \left\{ \{\mathbb{V}_1, \mathbb{N}_1\}, \{\mathbb{V}_2, \mathbb{N}_2\}, \dots, \{\mathbb{V}_n, \mathbb{N}_n\} \right\}$
7. Complete the $\{\mathbb{V}_j, \mathbb{N}_j\} \in \mathbb{V}$
8. Construct the \mathbb{V}^* set of sets
9. Reorder $\bar{\bar{\mathbf{L}}}$ in the a new, temporary matrix, $\bar{\bar{\mathbf{L}}}_0 \in \mathbb{R}^{m \times n}$.
10. Reorder the temporary matrix, $\bar{\bar{\mathbf{L}}}_0$ into the final matrix $\bar{\bar{\mathbf{L}}}^* \in \mathbb{R}^{m \times n}$.

The description of our algorithm for the case when $\bar{\bar{\mathbf{L}}}$ is composed of non-disjoint sets is given later in chapter 3.

8.1 Compute $\bar{\bar{\mathbf{A}}}_U$

We use the expression

$$\bar{\bar{\mathbf{A}}}_U = \bar{\bar{\mathbf{L}}} \bar{\bar{\mathbf{L}}}^\top \quad (1.4.1)$$

to compute $\bar{\bar{\mathbf{A}}}_U \in \mathbb{R}^{m \times m}$.

If a different weighting of the relations between the elements of $\bar{\bar{\mathbf{L}}}$ is desired, we can use the alternative expression

$$\left(\bar{\bar{\mathbf{A}}}_U\right)_{ij} = \sum_{k=1}^n \left\{ \delta^\star \left(\bar{\bar{\mathbf{L}}}, i, k \right) \bar{\bar{\mathbf{L}}}_{jk} \right\} \quad (1.4.1a)$$

where $i, j \in [1, m] \subset \mathbb{Z}^+$.

8.2 Construct \mathbb{U}

We initially construct \mathbb{U} by following the sub-routine (sub algorithm) given below

After the completion of this sub-routine, we will make a copy of this initial state of \mathbb{U} . We denote this copy of this initial state as \mathbb{U}_0 .

Data: Connection Matrix for rows, $\bar{\bar{\mathbf{A}}}_U$, and the number of rows in $\bar{\bar{\mathbf{L}}}$, m .

Result: Initial value for each of the $\{\mathbb{U}_i, \mathbb{M}_i\} \in \mathbb{U}$.

```

begin
  for  $i = 1 : m$  do
     $\mathbb{U}_i = \emptyset; \mathbb{M}_i = \emptyset$ ;
    for  $j = 1 : m$  do
      if  $(\bar{\bar{\mathbf{A}}}_U)_{ij} \neq 0$  then
         $\mathbb{U}_i = \mathbb{U}_i \cup \{j\}$ ;
         $\mathbb{M}_i = \mathbb{M}_i \cup \left\{ (\bar{\bar{\mathbf{A}}}_U)_{ij} \right\}$ ;
      end
    end
     $\mathbb{U}[i, 1] = \mathbb{U}_i; \mathbb{U}[i, 2] = \mathbb{M}_i$ ;
  end
end

```

Algorithm 5: Computing the initial value for each of the $\{\mathbb{U}_i, \mathbb{M}_i\} \in \mathbb{U}$.

8.3 Complete the \mathbb{U} Sets

The elements of \mathbb{U} must now be "*completed*" so that they include any indirectly related elements. Taking advantage of the disjointedness of $\bar{\bar{\mathbf{L}}}$, the "*completion*" of the elements in \mathbb{U} can be accomplished using the simple subroutine below

Data: Initial values for the elements of \mathbb{U} , $\{\mathbb{U}_i, \mathbb{M}_i\}$, and the number of rows in $\bar{\mathbb{L}}$, m .

Result: The final values for each of the $\{\mathbb{U}_i, \mathbb{M}_i\} \in \mathbb{U}$.

begin

```

    boolean isChanged = true ;
    while isChanged do
        isChanged = false ;
        for  $i = 1 : m$  do
             $\mathbb{U}_i = \mathbb{U}[i, 1]; \mathbb{M}_i = \mathbb{U}[i, 2]$  ;
            for  $j = 1 : m$  do
                 $\mathbb{U}_j = \mathbb{U}[j, 1]; \mathbb{M}_j = \mathbb{U}[j, 2]$  ;
                if  $\mathbb{U}_i \neq \mathbb{U}_j \ \&\& \ \mathbb{U}_i \cap \mathbb{U}_j \neq \emptyset$  then
                     $\mathbb{U}_i = \mathbb{U}_i \cup \{\mathbb{U}_j, \mathbb{U}_i\}$  ;
                     $\mathbb{M}_i = \mathbb{M}_i \cup \{\mathbb{M}_j, \mathbb{M}_i\}$  ;
                     $\mathbb{U}_j = \mathbb{U}_j \cup \{\mathbb{U}_i, \mathbb{U}_j\}$  ;
                     $\mathbb{M}_j = \mathbb{M}_j \cup \{\mathbb{M}_i, \mathbb{M}_j\}$  ;
                    isChanged = true ;
                end
                if isChanged then
                     $\mathbb{U}[j, 1] = \mathbb{U}_j; \mathbb{U}[j, 2] = \mathbb{M}_j$  ;
                end
            end
            if isChanged then
                 $\mathbb{U}[i, 1] = \mathbb{U}_i; \mathbb{U}[i, 2] = \mathbb{M}_i$  ;
            end
        end
    end

```

end

Algorithm 6: Computing the final value for each $\{\mathbb{U}_i, \mathbb{M}_i\} \in \mathbb{U}$.

8.4 Construct the \mathbb{U}^* Sets

To construct the sets of \mathbb{U}^* , we look for unique $\mathbb{U}_i \in \mathbb{U}$ and then store each unique \mathbb{U}_i as its own set in \mathbb{U}^* . We accomplish this by using the following subroutine

Data: The final values for the elements of \mathbb{U} , $\{\mathbb{U}_i, \mathbb{M}_i\}$, and the number of rows in $\bar{\bar{\mathbf{L}}}$, m .
Result: The sets of \mathbb{U}^* , with each set in \mathbb{U}^* representing a row cluster and containing its elements.

```

begin
  int nClust = 1;  $\mathbb{U}^*[nClust] = \mathbb{U}[1, 1]$  ;
  for  $i = 2 : m$  do
    boolean isNew = true;  $\mathbb{U}_i = \mathbb{U}[i, 1]$  ;
    for  $j = 1 : nClust$  do
       $\mathbb{U}_j^* = \mathbb{U}^*[j]$  ;
      if  $\mathbb{U}_j^* = \mathbb{U}_i$  then
        isNew = false ;
        break ;
      end
    end
    if isNew then
      nClust ++ ;
       $\mathbb{U}^*[nClust] = \mathbb{U}_i$  ;
    end
  end
end
end

```

Algorithm 7: Compute the set of all row clusters sets, \mathbb{U}^* .

8.5 Compute $\bar{\bar{\mathbf{A}}}_V$

We use the expression

$$\bar{\bar{\mathbf{A}}}_V = \bar{\bar{\mathbf{L}}}^\top \bar{\bar{\mathbf{L}}} \quad (1.4.2)$$

to compute $\bar{\bar{\mathbf{A}}}_V \in \mathbb{R}^{n \times n}$.

If a different weighting of the relations between the elements of $\bar{\bar{\mathbf{L}}}$ is desired, we can use the alternative expression

$$\left(\bar{\bar{\mathbf{A}}}_V\right)_{ij} = \sum_{k=1}^m \left\{ \delta^* \left(\bar{\bar{\mathbf{L}}}, k, i \right) \bar{\bar{\mathbf{L}}}_{kj} \right\} \quad (1.4.2a)$$

where $i, j \in [1, n] \subset \mathbb{Z}^+$.

8.6 Construct \mathbb{V}

We initially construct \mathbb{V} by following the sub-routine (sub algorithm) given below

Data: Connection Matrix for columns, $\bar{\bar{\mathbf{A}}}_V$, and the number of columns in $\bar{\bar{\mathbf{L}}}$, n .

Result: Initial value for each of the $\{\mathbb{V}_j, \mathbb{N}_j\} \in \mathbb{V}$.

begin

for $j = 1 : n$ **do**

$\mathbb{V}_j = \emptyset; \mathbb{N}_j = \emptyset;$

for $i = 1 : n$ **do**

if $\left(\bar{\bar{\mathbf{A}}}_V\right)_{ji} \neq 0$ **then**

$\mathbb{V}_j = \mathbb{V}_j \cup \{i\};$

$\mathbb{N}_j = \mathbb{N}_j \cup \left\{ \left(\bar{\bar{\mathbf{A}}}_V\right)_{ji} \right\};$

end

end

$\mathbb{V}[j, 1] = \mathbb{V}_j; \mathbb{V}[j, 2] = \mathbb{N}_j;$

end

end

Algorithm 8: Computing the initial value for each of the $\{\mathbb{V}_j, \mathbb{N}_j\} \in \mathbb{V}$.

After the completion of this sub-routine, we will make a copy of this initial state of \mathbb{V} . We denote this copy of this initial state as \mathbb{V}_0 .

8.7 Complete the \mathbb{V} Sets

The elements of \mathbb{V} must now be "completed" so that they include any indirectly related elements. Taking advantage of the disjointedness of $\bar{\bar{\mathbf{L}}}$, the "completion" of the elements in \mathbb{V} can be accomplished using the simple subroutine below

Data: Initial values for the elements of \mathbb{V} , $\{\mathbb{V}_j, \mathbb{N}_j\}$, and the number of columns in $\bar{\bar{\mathbf{L}}}$, n .

Result: The final values for each of the $\{\mathbb{V}_j, \mathbb{N}_j\} \in \mathbb{V}$.

```

begin
  boolean isChanged = true ;
  while isChanged do
    isChanged = false ;
    for j = 1 : n do
       $\mathbb{V}_j = \mathbb{V}[j, 1]$ ;  $\mathbb{N}_j = \mathbb{V}[j, 2]$  ;
      for i = 1 : n do
         $\mathbb{V}_i = \mathbb{V}[i, 1]$ ;  $\mathbb{N}_i = \mathbb{V}[i, 2]$  ;
        if  $\mathbb{V}_j \neq \mathbb{V}_i$  &&  $\mathbb{V}_j \cap \mathbb{V}_i \neq \emptyset$  then
           $\mathbb{V}_j = \mathbb{V}_j \cup \{\mathbb{V}_i\}$  ;
           $\mathbb{N}_j = \mathbb{N}_j \cup \{\mathbb{N}_i\}$  ;
           $\mathbb{V}_i = \mathbb{V}_i \cup \{\mathbb{V}_j\}$  ;
           $\mathbb{N}_i = \mathbb{N}_i \cup \{\mathbb{N}_j\}$  ;
          isChanged = true ;
        end
        if isChanged then
           $\mathbb{V}[i, 1] = \mathbb{V}_i$ ;  $\mathbb{V}[i, 2] = \mathbb{N}_i$  ;
        end
      end
      if isChanged then
         $\mathbb{V}[j, 1] = \mathbb{V}_j$ ;  $\mathbb{V}[j, 2] = \mathbb{N}_j$  ;
      end
    end
  end
end
end

```

Algorithm 9: Computing the final value for each $\{\mathbb{V}_j, \mathbb{N}_j\} \in \mathbb{V}$.

8.8 Construct the \mathbb{V}^* Sets

To construct the sets of \mathbb{V}^* , we look for unique $\mathbb{V}_i \in \mathbb{V}$ and then store each unique \mathbb{V}_i as its own set in \mathbb{V}^* . We accomplish this by using the following subroutine

Data: The final values for the elements of \mathbb{V} , $\{\mathbb{V}_i, \mathbb{M}_i\}$, and the number of columns in $\bar{\mathbf{L}}$, n .
Result: The sets of \mathbb{V}^* , with each set in \mathbb{V}^* representing a column cluster and containing its elements.

```

begin
  int nClust = 1;  $\mathbb{V}^*[nClust] = \mathbb{V}[1, 1]$  ;
  for  $i = 2 : n$  do
    boolean isNew = true;  $\mathbb{V}_i = \mathbb{V}[i, 1]$  ;
    for  $j = 1 : nClust$  do
       $\mathbb{V}_j^* = \mathbb{V}^*[j]$  ;
      if  $\mathbb{V}_j^* = \mathbb{V}_i$  then
        isNew = false ;
        break ;
      end
    end
    if isNew then
      nClust ++ ;
       $\mathbb{V}^*[nClust] = \mathbb{V}_i$  ;
    end
  end
end
end

```

Algorithm 10: Compute the set of all column clusters sets, \mathbb{V}^* .

Chapter 9

CbC Non-Disjoint Sets

Connection by Clustering on Non-Disjoint Algorithm info goes here

Part III

Time Complexity Analysis

Chapter 10

SVD Signs - Time Complexity

Intro for the Time Complexity Analysis of the SVD Signs Algorithm info goes here

Chapter 11

SVD Gaps - Time Complexity

Intro for the Time Complexity Analysis of the SVD Gaps Algorithm info goes here

Chapter 12

CbC Disjoint Sets - Time Complexity

Intro for the Time Complexity Analysis of the CbC Disjoint Sets Algorithm info goes here

Chapter 13

CbC Non-Disjoint Sets - Time Complexity

Intro for the Time Complexity Analysis of the CbC Non-Disjoint Sets Algorithm info goes here

Part IV

Space Complexity Analysis

Chapter 14

SVD Signs - Space Complexity

Intro for the Space Complexity Analysis of the SVD Signs Algorithm info goes here

Chapter 15

SVD Gaps - Space Complexity

Intro for the Space Complexity Analysis of the SVD Gaps Algorithm info goes here

Chapter 16

CbC Disjoint Sets - Space Complexity

Intro for the Space Complexity Analysis of the CbC Disjoint Sets Algorithm info goes here

Chapter 17

CbC Non-Disjoint Sets - Space Complexity

Intro for the Space Complexity Analysis of the CbC Non-Disjoint Sets Algorithm info goes [here](#)

Part V

Results, Comparisons, and Conclusions

Chapter 18

Comparison of Result Data

Intro for the results of the Data Comparison (comparing the initial data with the reordered data from each algorithm and the reordered data from each algorithm to the reordered data from the other algorithms) info goes here

Chapter 19

Comparison of Metrics on Data

Intro for the results of the Metrics on the initial datasets and their reordered counterparts from each algorithm (comparing results of the metrics on the initial data with the reordered data from each algorithm and the reordered data from each algorithm to the reordered data from the other algorithms) info goes here

Chapter 20

Comparison of CPU Run-Time

Intro for the results of the CPU time required on each dataset, for each algorithm
(comparison data taken from multiple runs of the same algorithm on the same dataset in
order to get a good statistical sample) info goes here

Chapter 21

Comparison of Memory Space Requirements

Intro for the results of the Memory Space required on each dataset, for each algorithm
(comparison data taken from multiple runs of the same algorithm on the same dataset in
order to get a good statistical sample) info goes here

Chapter 22

Conculsion

Intro for the Conclusions goes here