

Clustering for Graph Partitioning

Andrea Smith

May 3, 2018

About the Author

The author has earned her Bachelor's of Science degree in Computer Science from Missouri University of Science and Technology. She is working towards her doctorate degree in Computer Science, with a focus on graph data mining, from the same university, anticipating completing 2022.

Contents

| | | |
|----------|---|-----------|
| 1 | Executive Summary | 2 |
| 2 | Introduction | 3 |
| 2.1 | Graph Data Mining | 3 |
| 2.1.1 | Basic Definitions | 3 |
| 2.1.2 | Frequent Subgraph Mining | 4 |
| 2.2 | Graph Partitioning | 4 |
| 2.3 | Clustering | 4 |
| 3 | Project Specifications | 6 |
| 4 | Detailed Design | 7 |
| 4.1 | Class Objects | 7 |
| 4.2 | Algorithms | 8 |
| 4.2.1 | K-Means | 8 |
| 4.2.2 | Distance-K | 8 |
| 4.2.3 | Kernel Clustering | 8 |
| 5 | Experimental Results | 10 |
| 5.1 | Experiment Setup | 10 |
| 5.2 | Results Analysis | 10 |
| 5.3 | Difficulties | 10 |
| | Appendices | 12 |
| A | UML Diagram for Custom Classes | 12 |
| B | Pseudo-Code for the K-Means Algorithm | 13 |
| C | Pseudo-Code for the K-Means Cluster Initialization Algorithm | 14 |
| D | Pseudo-Code for the Kernel K-Means Refinement Algorithm | 15 |
| E | Pseudo-Code for the Distance-k Initialization Algorithm | 16 |
| F | Code | 17 |

1 Executive Summary

This project aims to explore the potential of clustering algorithms for aiding in the partitioning of large graphs for frequent subgraph mining. Partitioning is important as many graphs of interest cannot be held in the memory of a single machine, and so must be spread across many machines in order to be processed. Partitioning risks losing data that spans across multiple partitions, requiring an intelligent partitioning process in order to minimize.

This project explores k-means clustering, distance-k clustering, and multilevel kernel clustering algorithms, as they apply to partitioning large single graphs in data mining. The clusters produced by the algorithms are evaluated for quality based on cluster density and the maximal shortest path within the cluster. The algorithms were tested on a toy data set built from publicly available Arabian horse pedigree papers, such that each edge represents a child-parent relationship. The algorithms were written for undirected, weighted graphs, but could be modified for directed graphs.

Further research needs to be done to create and evaluate parallelized or distributed versions of these algorithms, so that they can process graphs that do not fit within the memory of a single machine.

2 Introduction

This paper is presenting basic research done to examine if clustering could assist in frequent subgraph mining for very large single graph settings. The size and density of many graphs of interest, particularly those involved with social media or chemistry, is a roadblock for efficiently dealing with them. Especially in situations where the graph is too large to fit in a single machine, we want to examine if clustering of the data can help make the process manageable.

2.1 Graph Data Mining

Graph data mining is a field of study concerned with extracting data from graphs. Most often, this data takes the form of a subgraph that has a structure that makes it interesting.

2.1.1 Basic Definitions

A graph is a set of nodes, V , and a set of edges, E . A graph can have labeled nodes or labeled edges, and both nodes and edges can also have numerical weights. Edges can also be directed edges, which mean that they come from one node and point to another node. They can be likened to a one-way road. For our purposes, we will only be concerning ourselves with undirected graphs.

A path in a graph is a sequence of nodes, such that every node has an edge that connects it to the following node. There can be many such paths between two nodes, but the distance between them is the number of edges that are in the shortest path between them. In a group of nodes, the longest distance between any two nodes (the longest shortest path) is referred to as the diameter of that group of nodes.

Graphs can be categorized as either connected or unconnected graph. If a graph is connected, that means that there exists a path to every node in the graph from every other node in the graph. If a graph is unconnected, then it is composed of multiple connected components, but there doesn't exist a path from one component to the other.

A graph also has a measurable density. Density is measured as

$$2 * |E| / (|V| * (|V| - 1))$$

for undirected graphs. The symbol $|x|$ means the number of items in that set.

So, $|E|$ would refer to the total number of edges. The number of edges connected to a single node, n , is known as degree and is also marked as $\deg(n)$.

2.1.2 Frequent Subgraph Mining

Frequent subgraph mining (FSM) is a common method for determining which structures in the graph occur most often. There are many different algorithms for doing this, which are out of scope for this project.

There are two overarching categories that a FSM algorithm may fall into. The algorithm is either meant to handle a single graph or a set of transaction graphs. For this project, we will only be concerned with single graphs, with particular interest in graphs which cannot fit within memory of one machine. While the intent of this project could be applied to a single graph from a transaction setting, that is not within the scope of this project.

2.2 Graph Partitioning

When dealing with a graph that is too large to fit within the memory of a single machine, some mechanism must be devised for dividing the graph among multiple machines. This process is known as graph partitioning.

Partitioning a graph is most frequently done by removing edges from the graph until there are multiple individual connected components. When doing this, it is important to remove as few edges as possible, but also make the partitions close to equal sized. One of the issues with doing a partition is that removing those edges makes it difficult to find frequent subgraphs that may span across multiple partitions.

2.3 Clustering

The goal of clustering is to divide a data set into logical groups of things that are similar. One option for defining similarity is distance, which would then allow us to apply clustering to large single graphs that would hopefully give as an approximation of how to partition the data

The three algorithms we will be testing in this paper are k-means, distance-k clustering, and kernel clustering. K-means is a traditional algorithm that works off of distance between elements [3], which we can define as the shortest path between two nodes. Distance-k is meant for graphs and clusters based on degree of the nodes [2]. Kernel clustering is a three step process where you make the data more coarse,

eventually allowing the graph to fit on one machine, applying a clustering algorithm, and then refining the data again back to the full size [1].

3 Project Specifications

This project focuses on comparing the quality of the clusters produced by each algorithm. Quality is a combined score of the cluster density and the maximum shortest path within the cluster. By comparing the clustering ability of the algorithms on a small graph, inferences about the quality of partitions made on a larger graph can be made. Note that further research should be done into parallelized and distributed versions of these algorithms, but that is beyond the scope of this project.

All code was written in the C++ programming language. The program accepts as input a properly formatted document file that describes the graph to execute. The program outputs the basic statistics of the clusters formed by each of the three algorithms, as well as a list of which nodes are in which cluster. The statistics reported include the number of clusters formed, the number of elements in each cluster, the density of each cluster, and the diameter of each cluster.

In the future, the program should be modified to accept a parameter file that can affect the execution of the program. Parameters of interest would include the number of clusters that the k-means algorithm should search for, and the sigma value for the kernel clustering algorithm.

4 Detailed Design

The program was designed in a combination of the object-oriented method and the functional method of programming.

4.1 Class Objects

The driving force of the program is a `UndirectedGraph` class, which is built from a `SymmetricalMatrix` class. The data is represented as an undirected graph because the project is looking for relationships between the horses, and we do not care about which horse is the parent. This then allows us to use a `SymmetricalMatrix` class, which is much more memory efficient to use than a full matrix class would be, as only either the upper or lower triangle of the matrix actually needs to be stored.

While C++ has a basic vector class, it does not have any linear algebra functions attached to that class. Instead, a custom vector class was developed for use, which the basic `Matrix` class was then built on top of. This combination allows for full flexibility to implement any of the needed linear algebra operations directly into the class. This was designed to be used by a spectral clustering algorithm, but that algorithm has instead been removed from the scope of this project. The `SymmetricalMatrix` class inherits from the `Matrix` class, but was given more efficient memory management, to take advantage of the fact the half of the values in the matrix are identical. A simple UML diagram for the relationship between the custom classes exists in Appendix A.

The `UndirectedGraph` class includes a `SymmetricalMatrix` object as a private member. This matrix is used to store the adjacency matrix of the undirected graph. The advantage to this method is realized in very dense graphs, as it becomes more efficient to store and to work with than other representations, which often involve long lists of node id pairs. Unfortunately, the graph used in the testing was not a particularly dense graph, and so this particular benefit did not come to light. However, dense graphs are a prioritized area of interest in the data mining field, and so it was important to consider that beyond this project.

The C++ data type `std::set` was chosen to represent the clusters within the functional algorithms. The class type has a strict enforcement of no duplication within a set and supports simple insertion, deletion, and search operators. Each node was represented in the cluster by its id in the original graph, not by its string label, which is more efficient as string operations are cumbersome and slow.

4.2 Algorithms

4.2.1 K-Means

The k-means function was broken into five different sub-functions. This helps improve readability of the code and allowed some sections to be reused for the other algorithms. The pseudo-code for the k-means algorithm is represented in Appendix B and the pseudo-code for the k-means centroid initialization is presented in Appendix C [3]. For this application, the distance between two nodes is taken as the shortest path between them, and the center of the cluster is selected to be the node with the highest closeness centrality. For future comparison, it would be interesting to see if different clusters were created if the center of the cluster was calculated based on betweenness centrality.

4.2.2 Distance-K

The distance-k algorithm is different from the other two, in that instead of taking a parameter that describes the number of clusters to create, it takes a parameter that describes the maximum size a cluster can be. This particular algorithm was developed for non-euclidean graphs, and so the size of a cluster is known as the diameter of it. This is more often known as the longest shortest path between any two nodes in the cluster.

The algorithm starts by forming clusters around the nodes that have the highest degree. It selects the highest degree node that is not already in a cluster, and forms a cluster of it and all of its neighbors who are not already in clusters. Pseudo-code for this process can be found in Appendix E. Once this is done, the algorithm will attempt to merge clusters together until they are too big, according to the distance parameter. Once no more merges can be made, the algorithm completes [2].

4.2.3 Kernel Clustering

Kernel clustering is done in three steps. The first step is coarsening the data. This step lends itself well to recursion, as each level of coarseness does not need to know about any previous or next level.

The second step is initial clustering. Since any clustering method can be done here, k-means was selected for this project. While k-means is reliable, it does potentially change the data quality. It would have been more complete to run the Kernel Clustering algorithm twice, once with k-means and once with distance-k, or to have selected a third clustering algorithm entirely for kernel clustering.

The third step of the algorithm is refinement. This requires a memory of the different levels of coarseness, in order to step back through them. Unfortunately, we did not think of a better way to solve this issue other than to hold in memory the different graphs created and a record of how the nodes were merged to reach that step. For large graphs, this could be a prohibitive amount of memory, especially with many iterations of coarsening.

At each level of refinement of nodes, the clusters are also refined. The algorithm presented for this was a modified k-means algorithm [3], which is expected to converge quickly at every iteration because the centroids received a good initialization due to the coarsening steps. The pseudo-code for the refinement of the clusters can be found in Appendix D.

5 Experimental Results

5.1 Experiment Setup

The experiment run was a graph with 64 nodes of Arabian horse pedigree relationships, where an edge between two nodes indicated a parent-child relationship. Due to the fact that only one of the tested algorithms had a random element to it, multiple iterations were not recorded.

5.2 Results Analysis

The results of each algorithm can be seen in Figure 1. Our standards for quality involve balancing the average cluster density and size against average cluster diameter. Ideally, the algorithms should find large and dense clusters, with a short path to all nodes in the cluster.

| | K-Means | Distance-K | Kernel Clustering |
|-----------------------|-------------|-------------|-------------------|
| Number of Clusters | 8 | 16 | 6 |
| Avg. Cluster Size | 8 | 4 | 10.66666667 |
| Avg. Cluster Density | 0.380254875 | 0.122179769 | 0.15488205 |
| Avg. Cluster Diameter | 3.75 | 1.25 | 4.333333333 |

Figure 1: Cluster Algorithm Results

It can be seen in the results that while Kernel Clustering returned the largest clusters, they also had the largest diameter, and a relatively low density. So rather than tight clusters, it found very loose clusters. Conversely, K-Means easily had the best density within it's clusters, and was on the high end of the cluster size. Again, however, it tended on the high end of the cluster diameters.

The Distance-k algorithm was notable in its very low average cluster diameter, but when examining the clusters themselves, the algorithm placed many nodes in clusters of one or two, thus very much lowering all three values collectively.

5.3 Difficulties

One of the main difficulties of this program was evaluating the calculated clusters for correctness. A visual aid tool was not developed, so all clusters had to be manually checked against the original matrix and drawn by hand during the development and

testing process. For continued work, a visualization tool would need to be integrated or developed.

Appendices

A UML Diagram for Custom Classes

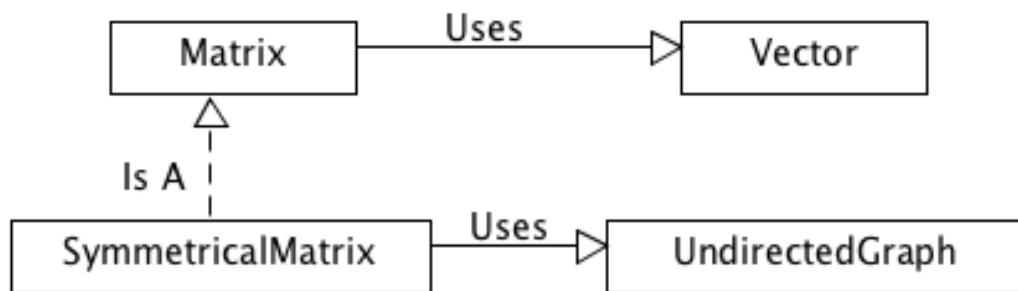


Figure 2: UML Diagram for Relationships between Class Types

B Pseudo-Code for the K-Means Algorithm

Input:
 $D = \{d_1, d_2, \dots, d_n\}$ // set of n data items

 k // Number of desired clusters

Output:

A set of k clusters.

Steps:

Phase 1: Determine the initial centroids of the clusters by using Algorithm 3.

Phase 2: Assign each data point to the appropriate clusters by using Algorithm 4.

Figure 3: Pseudo-Code for the K-Means Algorithm

C Pseudo-Code for the K-Means Cluster Initialization Algorithm

Input:
 $D = \{d_1, d_2, \dots, d_n\}$ // set of n data items
 k // Number of desired clusters
Output: A set of k initial centroids .
Steps:
1. Set $m = 1$;
2. Compute the distance between each data point and all other data- points in the set D ;
3. Find the closest pair of data points from the set D and form a data-point set A_m ($1 \leq m \leq k$) which contains these two data- points, Delete these two data points from the set D ;
4. Find the data point in D that is closest to the datapoint set A_m , Add it to A_m and delete it from D ;
5. Repeat step 4 until the number of data points in A_m reaches $0.75 \cdot (n/k)$;
6. If $m < k$, then $m = m+1$, find another pair of datapoints from D between which the distance is the shortest, form another data-point set A_m and delete them from D , Go to step 4;
7. For each data-point set A_m ($1 \leq m \leq k$) find the arithmetic mean of the vectors of data points in A_m , these means will be the initial centroids.

Figure 4: Pseudo-Code for the K-Means Cluster Initialization Algorithm

D Pseudo-Code for the Kernel K-Means Refinement Algorithm

WEIGHTED_KERNEL_KMEANS(K , k , w , t_{max} , $\{\pi_c^{(0)}\}_{c=1}^k$, $\{\pi_c\}_{c=1}^k$)

Input: K : kernel matrix, k : number of clusters, w : weights for each point, t_{max} : optional maximum number of iterations, $\{\pi_c^{(0)}\}_{c=1}^k$: optional initial clustering

Output: $\{\pi_c\}_{c=1}^k$: final clustering of the points

1. If no initial clustering is given, initialize the k clusters $\pi_1^{(0)}, \dots, \pi_k^{(0)}$ randomly. Set $t = 0$.
2. For each row i of K and every cluster c , compute

$$d(i, \mathbf{m}_c) = K_{ii} - \frac{2 \sum_{j \in \pi_c^{(t)}} w_j K_{ij}}{\sum_{j \in \pi_c^{(t)}} w_j} + \frac{\sum_{j, l \in \pi_c^{(t)}} w_j w_l K_{jl}}{(\sum_{j \in \pi_c^{(t)}} w_j)^2}.$$

3. Find $c^*(i) = \operatorname{argmin}_c d(i, \mathbf{m}_c)$, resolving ties arbitrarily. Compute the updated clusters as

$$\pi_c^{(t+1)} = \{i : c^*(i) = c\}.$$

4. If not converged or $t_{max} > t$, set $t = t + 1$ and go to Step 3; Otherwise, stop and output final clusters $\{\pi_c^{(t+1)}\}_{c=1}^k$.

Figure 5: Pseudo-Code for the Kernel K-Means Refinement Algorithm

E Pseudo-Code for the Distance-k Initialization Algorithm

```
begin
   $V' = V; P = 0;$ 
  while( $V' \neq \emptyset$ )
    begin
       $P = P + 1;$ 
       $m =$  the highest degree node in  $V'$ ;
       $Cluster(P) = \{m\};$ 
       $V' = V' - \{m\}$ 
       $\forall (m, n) \in E$  do
        begin
          if  $n \in V'$  then;
            begin
               $Cluster(P) = Cluster(P) \cup \{n\}$ 
               $V' = V' - \{n\};$ 
            end
          end
        end
      end
    end
  end
end
```

Figure 6: Pseudo-Code for the Distance-k Initialization Algorithm

F Code

The Git repository for this project can be found at:

<https://github.com/AriellaRomanov/Clustering>

References

- [1] I. Dhillon, Y. Guan, and B. Kulis, “A Fast Kernel-based Multilevel Algorithm for Graph Clustering,” *KDD*. Chicago, Ill.: 2005.
- [2] J. Edachery, A. Sen, and F. Brandenburg, “Graph Clustering Using Distance-k Cliques,” Arizona State University, 1999.
- [3] K. Nazeer and M. Sebastian, “Improving the Accuracy and Efficiency of the k-means Clustering Algorithm,” *Proceedings of the World Congress on Engineering*. London, U.K.: 2009.