Abdur Rafay Saleem
ab464825@ucf.edu

*Abstract*—**This review is based on a paper that proposes a new algorithm for data clustering that takes advantage of constraints. The structure of the paper is in three sections: the background and the research problem as well as the technical contributions, a detailed dive into the technical aspects of proposed algorithms, and a critical analysis of these techniques.**

I. RESEARCH PROBLEM

Clustering has many useful and interesting applications in the world of data and technology e.g. data analysis. A huge challenge grappling the data analysis field is the difficulty to cluster multi-dimensional datasets. Traditional clustering algorithms, such as K-Means, have been widely used due to their simplicity and efficiency. However, these methods often make assumptions that may not hold true in real-world scenarios, such as isotropic clusters. These algorithms simply apply the technique over the original feature vectors but in reality the results are far away from the true meaning of the data. This is due to the inclusion of unimportant correlations among many dimensions of these feature vectors. These can lead to suboptimal clustering results. Some research efforts in the form of dimension reduction and adaptive distancing techniques are used to project this data space into a low-dimensional feature space that is closer to its intrinsic meaning. The authors of the paper have identified this gap in the field and propose a novel approach called Locally Weighted Clustering.

The paper also identifies the recent trend in the field of utilizing partial information to aid in the unsupervised clustering process. This is primarily done using constraints and a strict set of policies to optimize the data partitioning steps by penalizing their violations. The authors extend the LWC algorithm to handle pairwise constraints, resulting in the Constrained Locally Weighted Clustering (CLWC) algorithm. This extension aims to utilize this partial information effectively. The paper improves the accuracy of the data-clustering process in two ways. Firstly, by assigning each cluster a local weighting vector that captures its unique structure. Secondly, it also incorporates constrained learning into the local weighting scheme by grouping data points from the constraint set and assigning each group to a cluster based on specific criteria.

The paper provides valuable technical contributions such as giving a solution to the multi-dimensional data clustering problem and providing optimizations. It provides a thorough theoretical analysis of their proposed methods which includes proofs for the optimal cluster centroids and weights. Also, the authors validate their proposed methods through extensive experiments on various datasets and reveal the advantages of their techniques.

II. TECHNIQUES

**Locally Weighted Clustering (LWC)**

The main idea behind LWC is to capture the local structures of data clusters in the feature space. Other clustering techniques such as K-Means, rely on global metrics for distance which are not that effective at capturing local structures. Therefore, this technique uses different multiple weighted distance metrics for the different clusters. Points in the cluster are used to determine a weighting vector for every cluster. These weighting vectors are used in re-scaling the distance from a data point to the cluster centroid. This adaptive distance metric is used to place data points to the nearest cluster. This technique is slightly different from Locally Adaptive Clustering (LAC) where the sum of weights is used as the constraint and the quality of its output greatly depends on determining the critical coefficient, which is hard. Whereas, LWC uses the product of weights between any cluster to be equal to 1 as the constraint. Therefore, LWC doesn't need user parameters or the regulation term for the weights. Optimal clustering is achieved by minimizing Sum of squared weighted distances between data points and their centroid using Lagrange Multipliers. This scheme ensures that the

weights are calculated efficiently and the centroid and local weighted constraint isn't affected by other clusters.

The dimensions are weighted based on the strength of correlation with the points. However, it is possible that in some cluster-dimension pair the value of the sum of squared differences is really small. To avoid troubles with weight computation, a threshold is used in case the value is very small or very large. The adaptive weight for each cluster can be extended by considering the inverse of their covariance matrix for Mahalanobis distance. However, some clusters may have a small number of points and this can cause a problem with accuracy of clustering. Therefore, to achieve stableness, the clusters are fitted to ellipsoid shapes.

The algorithm for this technique is iterative in nature and reduces the objective function in each iteration until convergence or max user specified iterations is reached. These are the steps:

1. Starts with a certain number of centroids and selects each of them using either *subset further first* or *Forgy initialization.* Upon selection, all weights are set to 1.
2. This step is about cluster assignments. Each point's local distance metric is used to place it in the closest cluster. This reduces the objective function with every new assignment.
3. This step recomputes and updates the centroids and weights of the cluster based on all the containing points.
4. Then the algorithm alternates between Step 2 (cluster assignments) and Step 3 (updating centroids and weights).

**Constrained Locally Weighted Clustering (CLWC)**

This is the primary technique proposed by the paper. Building upon LWC, the paper introduces Constrained Locally Weighted Clustering (CLWC) which integrates the local distance metric learning with constrained learning. The constraints are represented as instance-level constraints indicating whether the corresponding pairs of data points belong to the same cluster or not.

If two data points should exist in the same cluster in the output it is called a *Must-Link* constraint else a *Cannot-Link* constraint. This partial knowledge improves accuracy of clustering, especially semi-supervised clustering. This is a two stage process.

*1. Chunklet-based Assignment Strategy*

The main concept that allows the technique to respect constraints effectively is by using this assignment strategy. Basically, instead of assigning points to clusters one-by-one, it assigns them in bulk based on chunklets.

a. Every point starts with being a chunklet of size 1. Any data points with Must-Link constraints between them are merged in the same chunklet.
b. Any points with Cannot-Link constraints must exist in different chunklets. The remaining points which don't have constraints are put in an *isolated chunklet.*
c. In the end, the whole chunklets can be assigned to the cluster instead of individual points. For the isolated chunklet, points are assigned to clusters based on the minimum sum of squared distances between centroid and all points. This merging continues till all Must-Link constraints have been addressed.

This strategy reduces the chances of incorrect assignment and increases the probability of true output. The more data points a chunklet has, the more chance of being assigned to the correct cluster.

*2. Constrained Clustering*

Based on the chunklets produced by the first stage, the chunklet graph is constructed. Each chunklet acts as a vertex. Edge is added between vertices which have a cannot-link constraint in any

one of their points. This generated graph is used in the cluster assignment stage and implicitly affects the updates on weights and new centroids during iterations. The same algorithm above is followed except in step 2, it:

a. Examines memberships between all data points. Points without constraints are assigned to the closest clusters.
b. Then chunklets are picked in groups of 1 or 2 and assigned to clusters. This repeats until all chunklets are assigned. A max function is used in every assignment to pick the largest chunklet along with its largest unassigned neighbor. All the following steps make a joint decision for both chunklets.
c. For the assignment it reduces the search space by avoiding all clusters which have a neighbor assigned due to chances of violation of cannot-link constraints.
d. From the remaining clusters, the minimum sum of squared distances is calculated between their centroid and chunklets' points. The chunklets are assigned to the cluster where this value is minimum.
e. When no feasible assignment can be found for a chunklet or a pair of chunklets, a conflict occurs. The paper discusses how to handle these conflicts. In such cases, the algorithm tolerates some violations and assigns them to the closest clusters without considering the cannot-link constraints between them and their assigned neighbor chunklets.

The time complexity of this algorithm is better than conventional K-means and constrained K-means because of the fast convergence of the iterative algorithm. Moreover, the chunklet assignment strategy of CWLC produces a higher number of correct assignments and does not require user-specified parameters for fine tuning.

## III. CRITICAL ANALYSIS

The changing nature and increasing dimensions of data is making techniques such as data-clustering extremely difficult. This paper proposed a new technique to solve the major challenges that occur in clustering and to improve its accuracy beyond traditional means. The paper proves its efficiency and scalability using theorems.

The authors conduct several experiments using Rand Index and Normalized Mutual Information as the evaluation metrics to compare their technique's performance with other popular ones. In both unsupervised and supervised learning, this technique is superior to others based on fast convergence, clustering accuracy, lack of need for parameter tuning, low number of constraint violations and less performance degradation in case of small number of constraints.

Despite the promising results and innovative approach, the technique does have some limitations. The CLWC algorithm is still iterative in nature and could potentially lead to longer computation times with larger datasets. Its performance is also heavily dependent on the initial selection of centroids, which can lead to varying results. Moreover, the handling of constraint violations is not entirely clear and may lead to suboptimal clustering results like in the case of conflicts. The paper also does not address how the algorithm would perform with noisy data or outliers which are common in real-world datasets.

Further exploration could also be done on improving the efficiency of the algorithm by integrating it with other optimization techniques. Other approaches such as nonnegative matrix factorization and random walk techniques could potentially enhance the algorithm's performance and applicability. Additionally, more robust methods for handling constraint violations and conflicts could be developed to improve the accuracy of the clustering results.

In conclusion, the paper is well written and proposes an effective technique for solving data clustering challenges. The authors also focused on the most common challenges on each step and provided smart optimization and explanation to address these.