

Peer Review

GROUP 10 - algo_miners

SAMAY VARSHNEY - samay@iitg.ac.in - 180101097

PULKIT CHANGOIWALA - changoiw@iitg.ac.in - 180101093

SAI SUMANTH MADICHERLA - madicher@iitg.ac.in - 180101068

KOMATIREDDY SAI VIKYATH REDDY - komatire@iitg.ac.in - 180101036

CONTENTS

Contents

Peer Review 1

1

2

PEER REVIEW 1: GROUP 9

ROCK: A ROBUST CLUSTERING ALGORITHM FOR CATEGORICAL ATTRIBUTES

SHORT SUMMARY OF THE REPORT

Report starts with the introduction of the ROCK algorithm; then Section 2 describes the overview of static version of the algorithm; Section 3 covers the related work; Section 4 explains the incremental algorithm; Section 5 is about code structure; Section 6 and Section 7 presents data sets and evaluation results respectively.

'ROCK' is a clustering algorithm that uses the idea of linking data points instead of clustering them on the basis of distance between data points. The two key steps in the static algorithm are the computation of links and merging of clusters based on goodness measure. The incremental version of ROCK clustering algorithm adds new points to the original database in batches. Since the data set can be large, first a random sample is drawn from the data set. After that, the ROCK algorithm is run over the sampled points that employ links. Finally, the remaining data points on the disk are allocated to the relevant clusters.

Initially, each point is represented as a single cluster. For each of the cluster, i , we maintain a local heap which contains every other cluster j , such that $link[i, j]$ is non zero. The local heap is arranged in descending order on the basis of their goodness value. During each iteration, the pair of clusters with the highest goodness is merged and the goodness measure of the corresponding clusters is updated. The algorithm terminates either when the desired number of clusters has been reached or when all the clusters are well separated.

Incremental Algorithm: When a batch of b points is added to old k clusters we obtain total of $b + k$ clusters. For each of the b new cluster, we find cluster with the best goodness measure to merge, algorithm is repeated till there are again k clusters. After this, outliers are removed from the clusters.

Two real-life data sets Mushroom and Congressional Votes Data sets were used for testing ROCK and incremental ROCK, containing only categorical attributes. The metrics used for result comparison are running time, memory footprint, clustering quality and data distribution in clusters. For Mushroom data set, 1500 points were processed statically and 6624 were added incrementally while for Congressional Votes Data set, 350 were processed statically and 85 were added incrementally.

Different parameters like θ , k , n affects the clustering quality in different ways. On increasing θ , the number of links in the data set decreases and quality of links increases. For best clustering, k should be close to the number of the class labels. On increasing the number of points added, clustering quality decreases.

KEY STRENGTHS OF THE PROJECT

- The project explained why no existing incremental algorithm exists for ROCK even though the ROCK is a popular clustering algorithm and was presented in the year 2000.
- The proposed model results are highly similar to the original ROCK results with random sampling, to the point where the size of clusters is very similar. For example, in the Mushroom

dataset result, 17 out of 20 clusters have the same number of data points. Moreover, the size of the remaining clusters is also very similar to the ROCK algorithm.

- The proposed model explored the idea of batch insertion.
- The incremental model performs significantly better than static ROCK with high precision.

KEY WEAKNESSES OF THE PROJECT

- The incremental algorithm does not include deletion.
- There was a slight trade-off in time complexity for ease of implementation in the py-clustering library. The study implemented static ROCK based on the py-clustering library.
- It's accuracy and clustering quality decreases significantly on increasing the number of incrementally added points and size of the dataset.

SUGGESTED IMPROVEMENTS IN THE PROJECT

- The use of adjacency list instead of adjacency matrix can be explored to increase space efficiency.
- An improved similarity measure like the one suggested in QROCK can be used.
- Could have implemented Static ROCK parallelly (like DROCK) instead of sequentially to improve the time complexity.
- max_goodness could be calculated efficiently by using some other data structures like set or priority queue.
- Clusters are considered outliers if its size is less than or equal to five percent of the largest cluster. This value of five percent seems to have found through guess work. A better formulation or description of how one came up with this value would have improved the project.

SUGGESTED IMPROVEMENTS IN THE REPORT

- Report should contain a toy dataset over which we can verify that the algorithm runs correctly
- The flowchart can be made simpler by removing the use of C++ variables.
- Some of the urls can be made links which on clicking will direct to that url.
- Correctness of the algorithm should be included in the report.
- There should be some sources in the report (like email) through which one can contact the authors.
- Description of various parameters like $f(\theta)$, θ , m_a and m_b is missing from the report. Adding lucid description of the parameters aids in understanding the algorithm well.
- Clustering results are compared in tabular form by listing each clusters' description, this way seems inefficient and time consuming, a graphical way or some single value metric calculation would have been better.

OVERALL RATING

On analyzing the report of Group 9, we found the overall report and idea to be much satisfactory and worthy, hence for calculating overall rating we have divided the rating as follows.

- **Key strengths of the project:** We found strengths of the algorithm promising, the results were significantly better than static ROCK and time taken was comparatively smaller. Thus we are rating strengths 2.5 out of 2.5.
- **Key weaknesses of the project:** The main weakness of the project is that it does not contains methods for point deletion. As deletion is a crucial part of dynamic data set, thus

absence of it makes the algorithm incomplete. Therefore, we are rating algorithm 1.8 out of 2.5 in this aspect.

- **Suggested improvements in the project:** As a good project has very little scope of improvements but we found few very significant improvements which can be done in the project. Thus, we are rating it 2.2 out of 2.5.
- **Suggested improvements in the report:** Report has a lot of scope for improvements. Report could have been made more clear and lucid. Thus, we are rating it 2.0 out of 2.5.

Thus, overall rating of the project is $2.5 + 1.8 + 2.2 + 2.0 = 8.5/10$.