

# **Scalable K means ++ using sequential and parallel computing**

**-AbssZy**

## **Table of Content**

1. Abstract
2. Deliverables
3. Description
4. Explanation
5. Methodology
6. Implementation
7. Output
8. Conclusion
9. References

## **Abstract**

Our aim in this project is to implement the scalable K-means++ algorithm. This algorithm allows k-means clustering for larger number of data points efficiently. It is an efficient way to initialize centres Implementation of this paper. Now about the K-means algo which is one of the widely used data processing algorithms. However, the initialization in the algorithm is crucial for obtaining a good solution. K-means++ caters to this, by obtaining an initial set of centers that is close to the optimum solution. But k-means++ is a sequential algorithm and is inefficient over large data-set, k-passes need to be made over the entire data to find a set of good initial centres. It is important to reduce the number of passes made to get good initialization.

## **Deliverables**

At the end of this project, we have:

- Implementation of the algorithm in sequential format.
- Implementation of the algorithm in parallel format.

**Implementation:-** The implementation language will be Python. Core parts of the work will be done in Cython.

Unit tests will be added to validate against the existing clustering algorithms. Benchmark analysis will be done using the `make_blobs` dataset.

Documentation and examples will be added to the user guide and the API.

- The initialization function will be added in `cluster/k_means_.py`

## **Description**

Currently, k-means clustering algorithm does not perform very well on really large datasets. The scalable-kmeans++ algorithm aims to use a parallel implementation which makes choosing the centers faster. Also as more than one candidate centers can be chosen at a point, it makes this type of initialization less sensitive to outliers.

## **Explanation of Algorithm**

K means++ needs k number of passes over the data, for an application with large amount of data, k can value to 100 or even 1000. We want to pick each point independently so we do not want to choose points from a joint distribution we pick each point independently This is very useful for distributed implementation because each point independently gets to decide if it should be in the set of centres or not. This intuitively corresponds to updating your distribution much more infrequently which is a much more coarser sampling. It works as: when we have a distribution we pick some points from it and then update the distribution. More points are taken from the distribution we get after the above step. This is k means parallel a.k.a Scalable k-means++.

## Methodology

In this work we acquire a parallel variant of the k-means++ introduction calculation and exactly exhibit its viable .The primary thought is that as opposed to testing a solitary point in each go of the k-means++ calculation, we test  $O(k)$  focuses in each round and rehash the procedure for around  $O(\log n)$  rounds. Toward the finish of the calculation we are left with  $O(k \log n)$  focuses that shape an answer that is inside a steady factor far from the ideal. We at that point recluster these  $O(k \log n)$  focuses into  $k$  introductory places for the Lloyd's cycle. This introduction calculation, which we call k-means | | , is very basic and fits simple parallel usage.

### Lloyd iterations

In simple terms Lloyd iterations are nothing but randomly sampling  $k$  initial cluster centroids from our dataset. The process starts when some number  $k$  of point sites are initially placed in the input domain. These would be the vertices that are to be smoothed, as said above in some use cases they may be placed at random or by intersecting a triangular mesh which would be with the input domain.

## Implementation

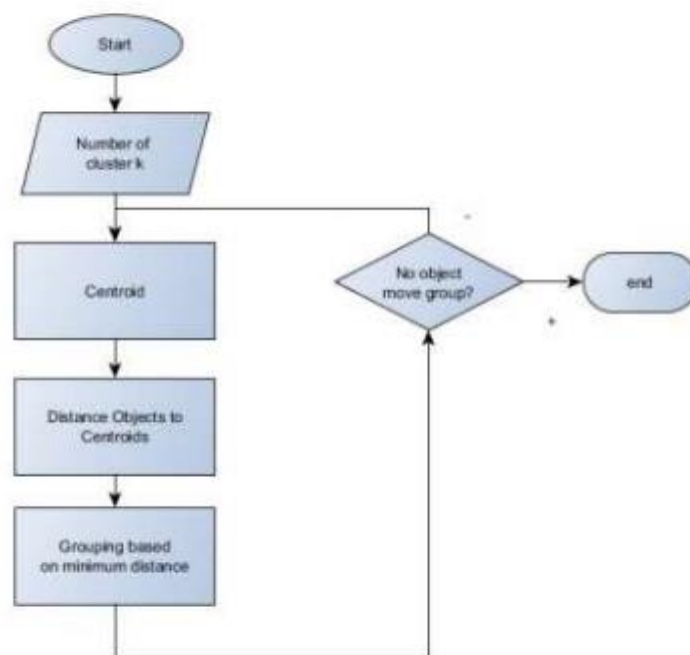
### Limitations

- Random initializations although fast would give rise to more number of iterations until convergence.
- Heavily susceptible to convergence at local minima

### K means++

K means is an algorithm used to cluster or gather together  $n$  number of objects based similar feature and attributes into  $k$  partitions. [ $k < n$ ] It has some aspects that are similar to the expectation-maximisation algorithm. The similar aspect is that both the algorithms try to find a natural cluster's centre in the data.

### How kmeans++:-



**Step A:** Let  $k$  = number of clusters

**Step B:** The partitions that classify the data are to be put into the k clusters. The initial centroid selection is done using a distance cumulative probability distribution function.

**Step C:** Get all the samples in a sequence and then calculate their distance from the centroid of each of their clusters. If the selected sample is not present in the cluster with the closest centroid, switch the sample to that cluster and update the centroid of the cluster.

**Step D:** Repeat step C until the centroids stop moving or convergence happens.

**Limitations:-**

- Asymptotic time taken is very long.
- This algorithm does not scale with data.



## **Conclusion**

As we can see in the results. Scalable Kmeans++ works much better than random initialization. Here the parallel version is computed with a relatively slow pace, but if more data is added it would work faster and more efficiently than the kmeans++ version.

## References

- [1] M. R. Ackermann, C. Lammersen, M. Martens, C. Raupach, C. Sohler, and K. Swierkot. StreamKM++: A clustering algorithm for data streams. In ALENEX, pages 173{187, 2010.
- [2] N. Ailon, R. Jaiswal, and C. Monteleoni. Streaming k-means approximation. In NIPS, pages 10{18, 2009.
- [3] D. Aloise, A. Deshpande, P. Hansen, and P. Popat. NP-hardness of Euclidean sum-of-squares clustering. Machine Learning, 75(2):245{248, 2009.
- [4] D. Arthur and S. Vassilvitskii. How slow is the k-means method? In SOCG, pages 144{153, 2006.
- [5] D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. In SODA, pages 1027{1035, 2007.
- [6] B. Bahmani, K. Chakrabarti, and D. Xin. Fast personalized PageRank on MapReduce. In SIGMOD, pages 973{984, 011.
- [7] B. Bahmani, R. Kumar, and S. Vassilvitskii. Densest subgraph in streaming and mapreduce. Proc. VLDB Endow., 5(5):454{465, 2012.
- [8] P. Berkhin. Survey of clustering data mining techniques. In J. Kogan, C. K. Nicholas, and M. Teboulle, editors, Grouping Multidimensional Data: Recent Advances in Clustering. Springer, 2006.
- [9] P. S. Bradley, U. M. Fayyad, and C. Reina. Scaling clustering algorithms to large databases. In KDD, pages 9{15, 1998}.
- [10] E. Chandra and V. P. Anuradha. A survey on clustering algorithms for data in spatial database management systems. International Journal of Computer Applications, 24(9):19{26, 2011}