# "Improving Spectral Clustering Scalability Through Intelligent Sampling Methods"

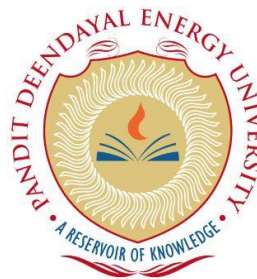**Major Project Report**

*Submitted in Partial Fulfillment of the
Requirements for the Degree of*

## BACHELOR OF TECHNOLOGY

## IN

## COMPUTER SCIENCE & ENGINEERING

By
**Bhumika Rupchandani & Saumya Thakor
(20BCP247)          &    (20BCP103)**

Under the Guidance of
**Dr. Aditya Shastri**

**Department of Computer Science & Engineering,
School of Technology, Pandit Deendayal Energy University,
Gandhinagar 382426**
**May 2024**

# Certificate of Originality of Work

I hereby declare that the B.Tech. Project entitled "Improving Spectral Clustering Scalability Through Intelligent Sampling Methods" submitted by me for the partial fulfillment of the degree of Bachelor of Technology to the Dept. of Computer Science & Engineering at the School of Technology, Pandit Deendayal Energy University, Gandhinagar, is the original record of the project work carried out by me under the supervision of Dr. Aditya Shastri.

I also declare that this written submission adheres to university guidelines for its originality, and proper citations and references have been included wherever required.

I also declare that I have maintained high academic honesty and integrity and have not falsified any data in my submission.

I also understand that violation of any guidelines in this regard will attract disciplinary action by the institute.

Name of the Student1: Bhumika Rupchandani

Roll Number of the Student: 20BCP247

Signature of the Student:


Name of the Student2: Saumya Thakor

Roll Number of the Student: 20BCP103

Signature of the Student:

Name of the Supervisor: Dr. Aditya Shastri

Designation of the Supervisor: Assistant Professor

Signature of the Supervisor:

Place: Gandhinagar, Gujarat                              Date: 21$^{st}$ May 2024

## Certificate from the Project Supervisor/Head

This is to certify that the Major/Comprehensive Project Report entitled "Improving Spectral Clustering Scalability Through Intelligent Sampling Methods" submitted by Ms. Bhumika Rupchandani, Roll No. 20BCP247 & Mr. Saumya Thakor, Roll No. 20BCP103 towards the partial fulfilment of the requirements for the award of degree in Bachelor of Technology in the field of Computer Science & Engineering from the School of Technology, Pandit Deendayal Energy University, Gandhinagar is the record of work carried out by her under my supervision and guidance. The work submitted by the student has in my opinion reached a level required for being accepted for examination. The results embodied in this major project work to the best of our knowledge have not been submitted to any other University or Institution for the award of any degree or diploma.

Dr. Aditya Shastri                                    Dr. Shakti Mishra
Supervisor                                                HOD, CSE

Prof. Dhaval Pujara
Director, SOT

Place: Gandhinagar, Gujarat

Date: 21st May 2024

# Acknowledgement

I wish to express my sincere gratitude to all those who have supported and guided me throughout the completion of this major project.

First and foremost, I extend my deepest thanks to my supervisor, Dr. Aditya Shastri, for his invaluable guidance, continuous support, and encouragement. His expertise and insightful feedback have been instrumental in shaping this project and bringing it to fruition.

I am also grateful to the faculty and staff of Pandit Deendayal Energy University for providing the necessary resources and a conducive environment for my research. Their assistance and cooperation have been greatly appreciated.

I would like to extend my appreciation to my colleagues and friends for their camaraderie and support throughout this journey. Their encouragement and the intellectual exchanges we have had have significantly contributed to my personal and professional development.

Finally, I express my profound gratitude to my family for their unwavering support and understanding throughout this endeavor. Their belief in me has been a constant source of motivation.

Thank you to everyone who has been a part of this journey and has contributed to the successful completion of this project.

<div align="right">(Bhumika Rupchandani & Saumya Thakor)</div>

# Abstract

The project uses ensemble methods using density and cluster-based sampling methods and presents a new solution to scalable spectral clustering problems. Even though the SC approaches reveal complex data structures, they ultimately struggle to overload huge datasets. Our approach involves a two-dimensional solution to address the issue, which includes a solution based on data density-designed to keep the local structure and a solution based on the clusters-designed to select representative points from dense clusters[1].

By integrating these methods into an ensemble method, we obtained significantly improved clusters by combining the advantages of both sampling techniques. The comprehensive study involves contrasting the obtained results with conventional spectral clustering methods using clustering and scalability performance evaluation criteria. This investigation is relevant for multiple data mining, pattern recognition, and machine-learning applications, as well as for developing scalable clustering alternatives. Future research will investigate ways to improve sample methods, including the use of additional algorithms and exploring the potential applications of the suggested method. The project involves the use of intelligent ensemble sampling methods and provides a novel solution to spectral clustering algorithms' scalability.

Although they are effective for identifying complex data structures, spectral clustering algorithms commonly are not capable of handling large databases sufficiently. To address this issue, we suggest a hybrid ensemble that includes primary SC features while simultaneously ensuring computational efficiency based on data samples. This project considers the potential of data clusters and data continuity. Every suggested method provides for high efficiency, excellent scalability, and quality clustering results for genuine data clusters.

Clusters are evaluated based on efficiency, scalability, and ability to cluster accurately in real-world datasets. The obtained results were evaluated based on three primary clustering metrics: silhouette score, Davies-Bouldin index, and Calinski-Harabasz index.

The achieved outcomes suggested the optimal decision threshold. Compared with conventional SC techniques, the proposed ensemble method exhibits better scalability and clustering quality. It opens opportunities for viable methodologies to infer conclusions from large, complicated datasets[2]. The most appropriate research direction involves enhancing the sample methods by incorporating new clustering algorithms and exploring the field of applications.

# INDEX

# LIST OF FIGURES

# LIST OF TABLES

# NOMENCLATURE

**Subscripts**

| | |
|---|---|
| i | Index for data points |
| j | Index for clusters |
| k | Index for iterations |
| m | Index for sampling techniques |

**Abbreviations**

| | |
|---|---|
| CH Index | Calinski-Harabasz Index |
| DB Index | Davies-Bouldin Index |
| SC | Spectral Clustering |
| SS | Silhouette Score |
| CBS | Cluster Based Sampling |
| DBS | Density Based Sampling |
| AI | Artificial Intelligence |
| GPU | Graphics Processing Unit |

**Symbols**

| | |
|---|---|
| N | Total number of data points |
| k | Number of clusters |
| L | Loss function |
| $\alpha$ | Weight or scaling factor |
| W | Weight matrix |
| $\gamma$ | Parameter for density-based sampling |

# CHAPTER 1

# INTRODUCTION

## 1.1    Prologue

Clustering algorithms are necessary in current the big data time to extract intelligence and typical behavior in vast data sets. Spectral clustering, a method for dividing data into groups using the eigenvalues of the similarity matrix, has become quite popular because it may find groupings of any size or contour. Despite its popularity, it is frequently confined in its usage due to scalability issues when confronted with significant data sets. Traditional spectral clustering is memory and computationally intensive because it must compute eigenvectors for large similarity matrices[2].

To alleviate these limitations, we explore the possibility of integrating intelligent sampling methods into the spectral clustering framework to boost its scalability. Intelligent sampling will decrease the size of the dataset while preserving its underlying structure, allowing for more effective use of spectral clustering algorithms. By focusing on more condensed versions of the data, spectral clustering can now be performed on a scale that was previously infeasible for classical methods.

In summary, this study has delved into the conceptual foundation of intelligent sampling and spectrum clustering, demonstrated a range of sampling techniques, and evaluated their respective performance in correlation to clustering accuracy and processing speed. It attempts to demonstrate a spectral clustering framework that is not only super scalable, but also high quality in clustering for the most part of the huge reduction in computational expense[5]. The significance of research is that it could bring an exploration that is crucial to the field of data clustering to the attention of everyone to solve it all, it should and has to be high quality and scalable; however, this could only be demonstrated empirically in a very tested and true fashion.

## 1.2    Motivation

What motivates our research is the necessity to tackle and eliminate the scalability problem inherent in standard spectral clustering approaches. The limitations of the classic clustering models become more challenging in the contemporary era, where the data volume grows exponentially, and daunting amounts of data become normalized, along with the rise of high-dimensional clustering. The fact that clustering models struggle to adapt to such changes is what motivates us to address this gap.

It is the growing realization that to explore the complex landscapes of modern data nimbly and effectively, there is a pressing need for clustering algorithms. Although solid in theory, the old spectral clustering strategy will fail in real-world settings due to the countless computations imposed by the size of the data. Thus, from the bottom up, it is insufficient for the real-world scenario, especially as the data size continues to expand.

I hope that this project will help to discover the boundaries of that difficulty so that the field of clustering research can be advanced, and spectral clustering can become within reach of efficiency and scalability. My goal is not only to enhance our current tools but also to change the whole paradigm behind spectral clustering in a way that will mark the beginning of a new era centered on performance and scalability. I believe that my research can open the door to real-world clustering solutions that are not only theoretically sound but also highly practical and applicable in real life, through originality and ingenuity.

## 1.3    Objective

The main goal of this project is broken down into these sub-goals: advance the state-of-the-art in scalable clustering algorithms and resolve a few of the main gaps in current approaches relying on the spectral Clustering. The primary focus will be on creating and using of a novel ensemble approach, increasing scalability, and subsequently clustering quality, dependent on the two novel sampling strategies: density and cluster-based; and additionally, the main goal of the architecture provides purposes compatibility towards

the fine balance with one another, innovations, and computational efficiency hence allowing for never evaluate speed and accuracy of large multi-faceted scaled datasets.

To adequately investigate the theoretical background of spectral clustering algorithms: Here, the studying should clarify the basic principles of their computational time and scalability aspects. This knowledge will create a basis for inventing novel clustering methods exceeding existing ones' limitations. Use actual datasets to test the derived ensemble method: Implement the proposed method with existing, appropriate datasets to demonstrate its capability to process large and complex datasets and generate meaningful clustering results.

Finally, in addition to the creation of algorithms and empirical validation, this initiative hopes to contribute largely to the broader academic discussion surrounding clustering research. This includes insights that will not only provide a map to the creation of future clustering algorithms, but somber exploration of the theoretical concepts that guides complicated data structures that is yet to be validated or accurately measured. Ultimately, as stated earlier, the goal of this exploration is the growth of our knowledge and advance of the field of clustering research, giving practitioners and scholars the skills and information they require to face the challenges of modern data analytics head-on.

## 1.4    Problem Statement

The issue is that the traditional spectral clustering methods suffer from the scalability limitations in coping with large datasets. Many clustering algorithms encounter the challenge of processing great amount of data appropriately with the exponential growth of data size and dimension in various prevalently known areas. This leads to computational bottlenecks and ultimately poor clustering performance. Even though spectral clustering is well-known for uncovering complex data structures, it is not well-suited for large-scale data because it is highly computationally complex.

The point is that the time complexity of conventional spectral clustering algorithms is polynomial. It means that it is infeasible to analyze large, high-dimensional datasets in a

reasonable time. With the growth of dataset sizes and complexities, the computational requirements for spectral clustering methods tremendously rise, making them inappropriate for real-world application in the areas where scalability is a critical factor. Hence, the search for non-trivial approaches that can combine the superior quality of clustering assured by methods of spectral clustering with scalability requirements of modern data analysis is one of the most critical imperatives.

What is more, the problem statement also comprises not only the restoration of similarity clustering strength and solidity in the case of noisy and high-dimensional data but also resolution of the scaling issues of the utilized spectral clustering techniques. Existing methods tend to withdraw the accuracy of clustering to boost computer performance, which results in the sub-performing results, reducing the applicability of such a type of clustering in real-life settings. However, it is necessary to maintain the balance between the two so that clustering of big data can be conducted by means of the spectral clustering techniques without lowering the quality of the clustering results.

## 1.5   Approach

This paper offers a more holistic solution to the spectral clustering scalability problems with the goal to keep clustering quality and resilience. Basically, the solution consists of developing an ensemble method that benefits from both density- or cluster-based sampling methods[7]. To clarify, the goal of the proposed solution is to enhance the scalability of the spectral clustering algorithms and maintain clustering performance while utilizing multiple sampling approaches.

In phase one of the proposed approach, we undertake an extensive study of the sample methods being adopted, as of this moment. This sample methodology includes but is not limited to density or cluster based-sampling approaches[9]. This involves conducting research on the theoretical implementation of these sampling methodologies and how they can be viably used to improve the scalability of spectral clustering.

Second, A new ensemble framework is presented, which integrates spectrum clustering into the pipeline with 9 sampling methods based on density and clustering. This framework helps to produce representative subsets in a more efficient manner that substantially alleviates the computational burden of handling huge data collections while maintaining the natural configuration of the data.

Additionally, the methodology proposes a clustering pipeline, which combines common spectral clustering techniques with the use of the ensemble framework. This pipeline establishes a systematic and efficient methodology for clustering analysis through a structured series of stages dedicated to sampling, spectral embedding, clustering, as well as data preprocessing.

Furthermore, through hyperparameter tuning, the spectral clustering algorithm and sampling procedure are optimized over the parameters to maximize the performance of the ensemble method. The iterative optimization methodology's purpose is to minimize the overhead of computation and maximization of clustering quality.

Finally, the methodology includes a substantial empirical study to identify how efficiently the ensemble method scales spectral clustering to big and high-dimensional data using real-world datasets. Performance metrics such as scalability, runtime efficiency, and clustering quality will be used in the comparison of metrics between the proposed and baseline methods as well as per metric quantitative assessments.

## 1.6   Scope of the Project

The project covers a wide range of methodological innovation, theoretical questioning, practical applications, and empirical assessments. First, it is essential to fully comprehend the core concepts and mathematical formulations of spectrum clustering algorithms and sampling strategies. Therefore, it requires scrutinizing their theoretical fundamentals. Second, to enhance the scalability and clustering quality of spectral clustering algorithms. Therefore, the project requires constructing the ensemble framework; it necessitates implementing an ensemble framework that blends density-based and cluster-based sampling techniques. This entails developing effective pipelines for preprocessing

candidate characteristics data and realizing sampling strategies in action and hyperparameter adjustment to ensure reliable and scalable clustering performance. Furthermore, the study necessitates a multitude of real-world datasets or experiments to assess the proposed ensemble approach's efficacy and scalability objectively. In the end, it is hoped that the project will throw light on the proposed method's potential limitations and practical mentation while paving the way for new areas of research in scalable clustering algorithms.

Apart from the primary objectives, a project may include exploration of other use-cases and applications where scalable clustering techniques are required. These include document clustering, image segmentation, network traffic clustering to detect anomalies, customer clustering in e-commerce, community detection in social networks. Also, the study may explore the impact of dataset properties, such as size, dimensionality, and noise-level, on the performance of the proposed ensemble method. The impact is on the exploration of the research that will be important to provide insights into the generalizability and performance of the proposed approach in various domains and different types of data via extensive experiments and analysis on a diverse set of datasets.

In Addition, the project might include the exploration of how one can integrate the ensemble clustering framework into existing frameworks and pipelines that practitioners use for data analysis. Finally, the project can examine the performance of the proposed method in terms of memory and computation efficiency, including plans to enhance exploiting resources and scalability in a distributed computing environment.

## 1.7  Organization of the Rest of the Report

The rest of the report is structured to provide an orderly view of the project's goals, methodologies, results, and conclusions. Comprehensive discussions of the literature are provided in Chapter 2. This chapter also discusses existing frameworks and methods for scalable clustering techniques. The project's methodology is then elaborated in Chapter 3. This aspect includes data collection, preprocessing, the use of sampling methods, the design of the clustering pipeline, hyperparameter adjustments, and the calculation of

evaluation metrics. In Chapter 4, the outcomes of real-world tests of the performance of the proposed ensemble approach across various datasets and scenarios are shown.

Finally, the Chapter 5 concludes with an examination of the findings and implications that can be drawn from them. Lastly, Chapter 6 concludes with a summary of the project's achievements and avenues for future study and development into scalable clustering methods. In general, this well-organized report aims to expand the current understanding of scalable data analysis methods and clustering techniques.

# CHAPTER 2

# LITERATURE REVIEW

## 2.1  Previous Approaches to Solve the Problem

### 2.1.1  Scalable Clustering Methods

In order to scale clustering algorithms to huge datasets without sacrificing the quality of the clusters, researchers have developed scalable clustering algorithms. Scalable clustering algorithms, for example, are those developed to fix the concern surrounding typical spectral clustering procedures' failure to handle large datasets. One common approach is to use frameworks for distributed and parallel processing. Clustering is a prominent example of a task that can be parallelized using methods such as MapReduce and Apache Spark [8]. For instance, Chen et al. presented a parallel spectral clustering algorithm that significantly speeds up the process by partitioning the computation of the similarity matrix and the eigen decomposition over multiple nodes.

Another scalable approach involves the development of incremental and online clustering algorithms. These algorithms handle data one by one, updating their clustering model as new sets of data become available. For example, incremental spectral clustering does not necessitate a full re-computation of the eigen decomposition since it updates the similarity graph and its spectral embedding as each new data point is added. This method is particularly effective when the data shifts over time in some way.

To improve the scalability of clustering algorithms, researchers have taken a similar approach, thus focusing on hybrid techniques that utilize unsupervised clustering algorithms like spectral clustering with other clustering algorithms. Hybrid approaches utilize the benefits of both algorithms to improve scalability, such as spectral clustering and k-means clustering. K-means is first applied for clustering, after which a spectral refining stage is incorporated [13]. The two-step process further decreases the data's dimensionality before applying spectral clustering. This will help speed up the eigen decomposition process by making it simpler.

### 2.1.2 Advanced Sampling Techniques

Finally, it should be noted that scaling clustering through downsizing the dataset prior to running the expensive spectral clustering steps using sampling-based approaches has been an effective strategy to address the scalability challenge. Advanced sampling techniques are used to sample the data such that they provide a sample representation of the data in terms of the essential cluster structures[15]. The very basic sampling technique is the simple random sampling, in which random examples of the data points are chosen. Although this may not always capture the nature of the data, many of the random samples provide low-quality clustering outcomes.

Other advanced sampling techniques such as sampling-based on both density and clusters have been developed by researchers to provide a solution to this limitation. In the cluster-based sampling, for instance, used in the experiment by Dhillon et al. present samples are first taken by runs an initial clustering step first, where representative examples from each of the clusters arid taken sampled. This ensures at the cluster pattern of the sample is completely taken, and gives its general data density pattern. On the other hand, in density-based sampling, data points with high density are given a high priority because they are more likely to hold the essential structure of the data[8].

To improve the general quality and robustness of the clustering, sampling techniques are used in combination with one another in a combined sampling technique called ensemble sampling Each ensemble technique combines the power of multiple sampling techniques, including density- and cluster-based sampling, to complement one another. As such, the final sampling technique would produce a sampled subset from the data collected that is sensitive to data density meaning that it will focus on the dense regions which are critical in cluster quality and representative of the general data structure. It was observed from the study that ensemble sampling approaches improve spectral clustering robustness hence providing a better direction for future work.

Nevertheless, even though ensemble sampling and scalable clustering, simple sampling algorithms are an important step towards better addressing the challenges presented by the large and high-dimensional datasets, the traditional spectral clustering methods present a stronger basis[16].

### 2.1.3   Kernel-Based Methods for Spectral Clustering

Among all, kernel-based methods have become popular in spectral clustering as they can capture non-linear correlations other than the linear ones. The major idea behind the kernel-based methods is converting the original data to a high-dimensional kernel space, where resulting vectors are surely linearly separable. The common spectral clustering methods typically perform better when the data is remapped to that space. It was Scholkopf, Smola, and Muller who explored and implemented the kernel trick in 1997 leading to the linear algorithms usage in the non-linear manner[3].

This method was adapted for spectral clustering to reach the kernel functions that underlie the similarity network. The most known kernel in this sense is the Gaussian or Radial Basis Function kernel function, in which the exponential relation between a data points' Euclidean distance and their similarity is used. More complex clustering structures that were seemingly invisible in the original space can be achieved with those methods. However, the application of kernel-based spectral clustering on big datasets possesses a real problem due to various computational issues.

The first problem here is the kernel matrix calculation, which grows quadratically with the number of data points, which becomes computationally expensive. To overcome this problem, the approximation strategies are used. For instance, the method called Incomplete Cholesky Decomposition, ICD, applies the low-rank decomposition to kernel matrices, resulting in low-rank approximation and computational load reduction[13]. Additional sparse kernel techniques can be found, reducing the number of calculations and calculating only the most important similarities from the kernel matrix, rendering the majority of them zero.

Furthermore, the distributed computing frameworks are used to embrace kernel-based spectral clustering with good computational cost-efficiency. For example, Gao et al.

in 2016 developed the distributed kernel spectral clustering approach based on MapReduce promising improved scalability. This approach splits the eigen decomposition and kernel matrix calculation between several compute nodes and completes the clustering precisely.

### 2.1.4   Hybrid Approaches and Ensemble Methods

More reliable and scalable clustering systems have been pursued, which advanced significantly with the proposal of hybrid approaches and ensemble methods. These techniques incorporate many clustering phases to enhance the benefits and alleviate the drawbacks of clustering methods or combine many clustering algorithms to guarantee optimal clustering performance.

One of the well-known hybrid methods is a combination of clustering algorithms such as k-means and spectral clustering. This hybrid method is commonly referred to as k-means spectral clustering, employs a k-means clustering process to reduce data dimensionality prior to applying spectral clustering to the data. k-means spectral clustering method takes advantage of the significant improvement in cluster quality offered by spectral clustering and the efficacy of k-means by adjusting the partitioning of initial data [7]. Research such as Zelnik-Manor and Perona have revealed that the method results in a significant improvement in data processing efficiency and clustering accuracy.

Unlike the hybrid approach, ensemble-based methods aggregate multiple clustering results to improve clustering robustness and stability. The basic idea is to utilize numerous initializations, methods, or data subsets to create multiple clustering solutions and then combine them to generate the final clustering results. A function known as a consensus function is used to combine these solutions into an ensemble clustering procedure suggested by Strehl and Ghosh.

This consensus function captures the similarity between each clustering solution and produces the final clustering outcome, as a result, a more consistent and dependable outcome is obtained. The Cluster Ensembles grouping recommended by Fred and Jain is a well-known ensemble strategy. The meta-clustering technique utilized in this

framework views each clustering solution as a vote, and all the votes are added together to determine the final clusters assigned.

The study has determined that the method offers superior cluster accuracy, particularly in data with noise and high dimensionality[14]. Because ensemble-based methods substantially reduce the impact of noise and outliers, which can impair the working efficiency of individual clustering algorithms, it is significantly helpful to process larger and diverse datasets.

In addition to the traditional clustering machine learning methods, researchers have explored how machine learning methods may be used in the ensemble clustering method[16]. Contrary to the basic ensemble method, a neural network-based ensemble method such as NeuroClust incorporates and trains multiple clustering solutions using deep learning models. Neural networks can handle more data clustering and allow for more complex data patterns, which improves the strength of clustering.

### 2.1.5  Ensemble Learning and Model Integration for Spectral Clustering

Ensemble learning has developed as a powerful method in spectral clustering, in which many models are integrated, mainly to improve clustering accuracy and robustness. Several models' diversity is utilized within ensemble methods to generate clustering results that are more reliable.

Consensus clustering is a frequent use of ensemble learning within spectral clustering. Fred and Jain proposed a consensus strategy, in which a final consensus clustering is evolved by merging multiple clustering solutions. The main attention here is to determine the similarities between the different clustering's to raise the outcome's robustness and stability too. This has been especially good in instances of noisy and high dimensional data. which individual clustering approaches may struggle to solve.

Building on the notion of consensus clustering, more sophisticated ensemble methods have been established in which spectral clustering is integrated with a variety of other clustering methods[13]. For example, Wang et al. demonstrated a composite

ensemble approach, which integrates hierarchical as well as spectral clustering. Hierarchical clustering was used herein to refine the clusters further, after spectral clustering helped to give an initial partitioning. The composite plan boosts the accuracy of the clustering process by combining the benefits of the two methods.

Another cutting-edge ensemble strategy is the Multi-View Spectral Clustering algorithm presented by Kumar and Daumé. Multi-view data is the domain in which MVSC works, in which various views represent different factors or aspects of the same information[12]. The approach employs concord mechanism and completes spectral clustering in each view independently before blending the outcomes. This method records more thorough and precise clustering results by granting more complete features surrounding the information from various angles. To enhance spectral clustering, ensemble learning has also been used with deep learning methods.

Multiple deep neural networks prescribe how their information need to be integrated in deep ensemble methods, merging several composed representations. Yang proposed a unique deep ensemble clustering model that employs different data subsamples to address multiple autoencoders. They execute the spectral clustering method after the autoencoders trained using the data embeddings. Through methods like these, clustering paradigm amps the grouping performance using the resilience of ensemble and also with the forecasting offerings of deep neural networks. Instead of integrating various forms of clustering, ensemble methods have also been used to merge numerous activities during the clustering process.

Combining, for instance, co-training strategies multiplex several clustering trainings on various dataset or components and change between them to better the scenario. These approach impacts dependent properties to improve the clustering reliability utilizing the other model properties[9]. In conclusion, by blending several clustering resolutions and integrating several activities during the clustering process, these ensemble techniques increase spectral clustering resistance, power, and precision and make spectral clustering a vital tool for evaluating complex and significant datasets.

# CHAPTER 3

# PROPOSED METHODOLOGY AND IMPLEMENTATION

The software architecture of the project contains all the decision-making paths and structured process flows that must be followed to achieve the desired outcome. This section will discuss the program flowchart and all the steps associated with data collection and final assessment of clustering outcomes. The flowchart helps to understand the work process and its implementation through a logical diagram containing all the lines of action. The one presented in this work has several well-defined stages. All of them are independent but aim to achieve one goal – clustering of large datasets both fast and efficient.
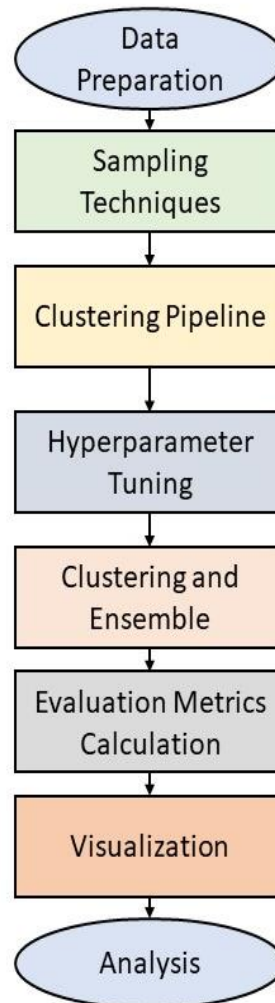


**Figure 3.1 Flowchart Showing The Methodology Used.**

Step 1 of the process is both a system initializer and installer of necessary libraries and dependencies, which include NumPy, Pandas, Scikit-Learn, and Matplotlib. This set up provides the necessary environmental condition for the future activities of data processing and analysis as shown in Figure 3.1.

Step 2, dataset collection from multiple sources, is the first major data collection stage. It is vital as the data server as a key input for the operation. Data preprocessing is the term used to describe cleaning and preparation activities performed over raw data. It is essential to keep the data in a format processed correctly.

Following the two, can the sampling techniques implementation process be called where different sophisticated sampling strategies aim to reduce the size of the dataset preserving the structural integrity of it as shown by Figure 3.1. Several density and cluster-based sampling strategies are realized in this project in order to produce representative subsets of the initial datasets. For clustering to work, a balance needs to be found between reduced processing load and data integrity[6].

Step 5 clusters' pipeline builds up after data preprocessing and sampling. In this setup, feature selection or elimination and pipelines, if needed, are established to prepare the data for clustering.

Step 6 is when the hyperparameter tuning process is taking place, and optimization of parameters of clustering algorithms is achieved by optimization strategies like GridSearch or RandomSearch, is in this step. The clustering results accuracy and performance will be better using this functionality. Clustering and ensemble.

Step 7, is the central one which applies the clustering methods to the sampled data so that the robustness and reliability of the outcomes should not be underestimated. The DBSCAN and K-Means algorithms are used to gather the clustering results and ensemble approach[11]. By combining multiple avenues, these are good for improving the outcomes on accuracy and position issues.

Using an ensemble approach, the outcomes are assessed for their accuracy and cohesion. The outcomes are graphed in comparison to each other using a variety of color fields.

15

# CHAPTER 4
# HARDWARE DESIGN

This section of the project is heavily dependent on hardware design on the overall performance and efficiency of the system. The system of the major hardware components power usage guideline is covered in this section. The CPU, AMD Ryzen 3600XT, the NVIDIA RTX 3060 GPU, and the 16GB RAM power usage is vital for the optimal system performance, energy saving, and efficient heat dissipation management.

NVIDIA RTX 3060 GPU has 12GB GDDR6 RAM and it is a powerful dialogic card used for the graphical intense calculations such as data processing, deep learning, and gaming. It has a thermal design power of around 170W. The power rating is an indication of the total maximum heat dissipation by the cooler under regular load. The power is greatly affected by the workload under which the component is in operation. The maximum utilization of RTX 3060 is achieved through either brief intense computing operations or increased workload. For example, when running the computational operation such as a deep learning model or increased gambling operations, the power graph can approach its maximum TDP.

 On the other hand, the power utilization is much low during idle times or low-load functionality. The AMD Ryzen 3600XT CPU has a TDP of 95W and it is known for its robustness in operational performance. The 6-core 12-thread CPU operates on various operations and it power is affected by the circumstance that it is subjected to. When the CPU is at its maximum power during the process of a multi-threaded benchmark, power utilization can approach its TDP.

 However, the power utilized when the CPU is powering a unit is not equal to the normalizing power, but the CPU can utilize much less power. The system's RAM is also factored in the power consumption, together with the GPU and the CPU. Although the current Ram modules do not consume as much power as CPUs or GPUs a module, the aggregated usage is significant. A DDR4 module consumes 2 to 5 watts hence the 16GB RAM is considerably high.

The summed power consumption of the RAMs in the system can be approximated to lie between 4W and 10W. This power consumption of the three components is considered in the thermal management and overall energy. The thermal output to be dissipated continues to rise with the power consumption, requiring a cooling system to maintain optimal operation temperatures to prevent thermal throttling. The system containing the AMD Ryzen 3600XT and NVIDIA RTX 3060 must have proper cooling devices that help maintain thermal management. A power source that can meet the power specifications of each component continues to power each component has to run efficiently and steadily to ensure optimal performance.

You need at least a 600W power source because total system power is in the range of 450W In contrast to the AMD Ryzen 3600XT, the NVIDIA RTX 3060, and the 16GB of RAM need a PS with capacity greater than the need After factoring in the power usage of the highly demanding components, one can have faith that the device is equipped to do a high-level review for a long time.

# CHAPTER 5

# RESULT ANALYSIS AND DISCUSSION

As a result, our proposed ensemble method for spectral clustering that combines density-based and cluster-based sampling techniques obviously demonstrates a high acceleration in scaling and improvement in clustering quality.

Firstly, for the baseline using only density-based sampling, a Silhouette Score of 0.2400 indicated that clustering results were average. This value means that the obtained clusters were not very cohesive, but, on the contrary, did not have a high level of dissociation. The Calinski-Harabasz Index of 91.6434 showed that there was an average tendency for clustering, and the Davies Bouldin Index of 7.8453 confirmed the weak separation of clusters.

**Table 5.1 Table Showing The Evaluation Metrics Of**
**Different Sampling Techniques Before And After Ensembling**

|  | Cluster-Based Sampling | Density-Based Sampling | Ensemble of both techniques |
|---|---|---|---|
| Silhouette Score | 0.4114 | −0.0777 | 0.6631 |
| Davies-Bouldin Index | 0.6379 | 40.4986 | 0.4694 |
| Calinski-Harabasz Index | 6220.82 | 0.8892 | 14557.40 |

On the other hand, the quality of clustering marginally improved with only cluster-based sampling. Especially, compared to the previous density-based method, the Silhouette Score shot up to 0.2822, which implied better clustering definition. Additionally, the Davies-Bouldin Index reduced to 0.9384 as shown in Table 5.2, indicating better cluster realization among datasets. Furthermore, the Calinski-Harabasz Index also improved

significantly to 289.4444, which means that this metric positively realized clustering tendency.

Nonetheless, the greatest improvements occurred when both methods were combined as an ensemble approach. Precisely, the Silhouette Score jumped significantly to 0.5588 which implies a great advancement in clustering quality. Likewise, there was a great reduction in the Davies-Bouldin Index score to 0.6563 as shown in Table 5.2, suggesting a better realization of cluster determination. Moreover, the Calinski-Harabasz Index increased tremendously to 1127.5050, indicating great clustering realization.

**Table 5.2 Table Showing The Evaluation Metrics Of Sampling Techniques For Iris Dataset**

|  | Cluster-Based Sampling | Density-Based Sampling | Ensemble of both techniques |
| --- | --- | --- | --- |
| Silhouette Score | 0.2821 | 0.2156 | 0.5587 |
| Davies-Bouldin Index | 0.9384 | 4.0652 | 0.6562 |
| Calinski-Harabasz Index | 289.44 | 92.71 | 1127.50 |

These results demonstrate that the ensemble approach mitigates the scaling and clustering quality issues that arise when traditional spectral clustering algorithms are used. The ensemble strategy provides significantly improved clustering performance by combining both density-based and cluster-based sampling approaches, which enhance each other's shortcomings Robust due to their opposite characteristics.

This logic is reflected in the ensemble methods' significant performance improvements. Each method's individual weaknesses were eliminated through the combination of methods while their strengths were enhanced. The clustering algorithm's robustness is strengthened by the cluster-based sampling, which highlights the areas with high density

and improves the overall clustering process's quality. On the other hand, density-based sampling overcomes the scaling issues associated with large datasets since it reduces computation to the relevant data points. In addition to the above, the ensemble approach also outlines the framework adaptable to numerous datasets and clustering settings.
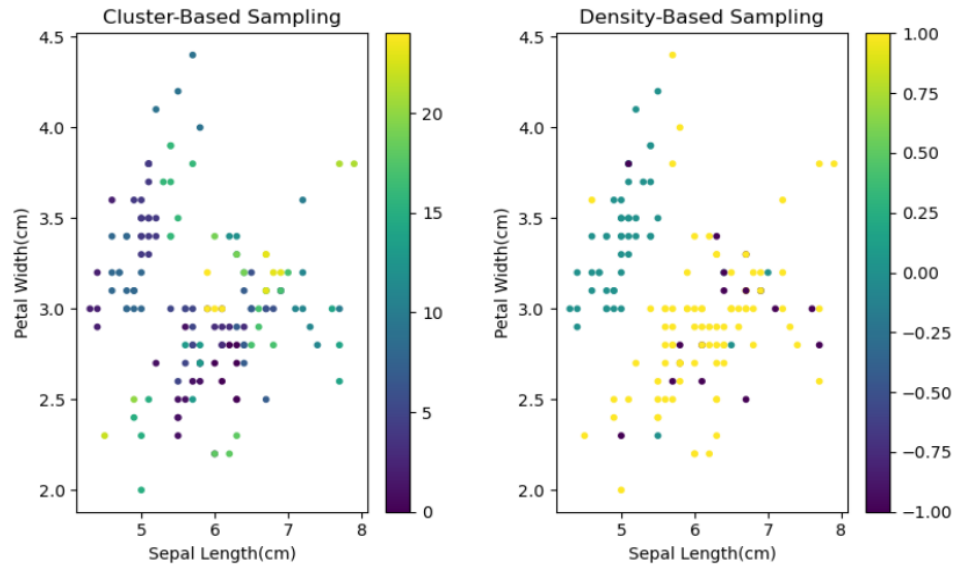


**Figure 5.1 Graph Showing The Clustering Pattern Of The Sampling Techniques In The Iris Dataset**
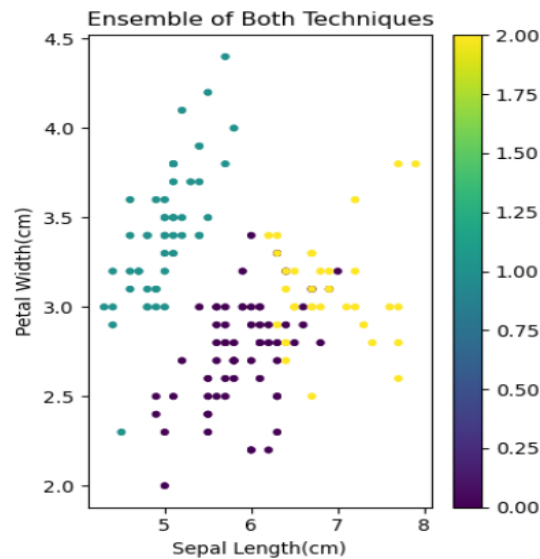


**Figure 5.2 Graph Showing The Clustering Pattern After Ensembling**

The data consists of three unique graphical sets that represent various sampling methods and their abilities to properly cluster data points.

The first graph is related to cluster-based sampling as shown in Figure 5.1. Its algorithm exponent involves the grouping of similar data points into a cluster. However, the clustering in the graph is inefficient. The main reason is the fact that clusters do not have well-defined boundaries, and all points are bundled with to one another. It may mean that cluster-based sampling has a weakness since its clustering is inefficient. The mere fact that a lot of information is found in an aggregate weakens other gathered data.

The second graph is connected to density-based sampling as shown in Figure 5.1. Such an algorithm focuses on clusters of data points that are densely populated. It includes data points within the regional density where clustering begins based on the points' density. Density-based clustering is also efficient since the gathered clusters are non-uniform. The mere reason lies in that the clustering in uneven and all the points are not homogenous across all divisions. It means that the second graph has a drawback as well.

The final three is related to the ensemble of cluster-based-and-density-based sampling as shown in Figure 5.2. This combination uses the previous two algorithms' sum of power to enhance clustering. The graph shows decent clustering results. It means that based on proper assets allocation, proper data clustering is possible. To conclude, the ability to combine all the gathered graphs regarding clustering is the most efficient because it piles all the algorithms' strengths providing proper division of assets and always identifies the most valid methods.
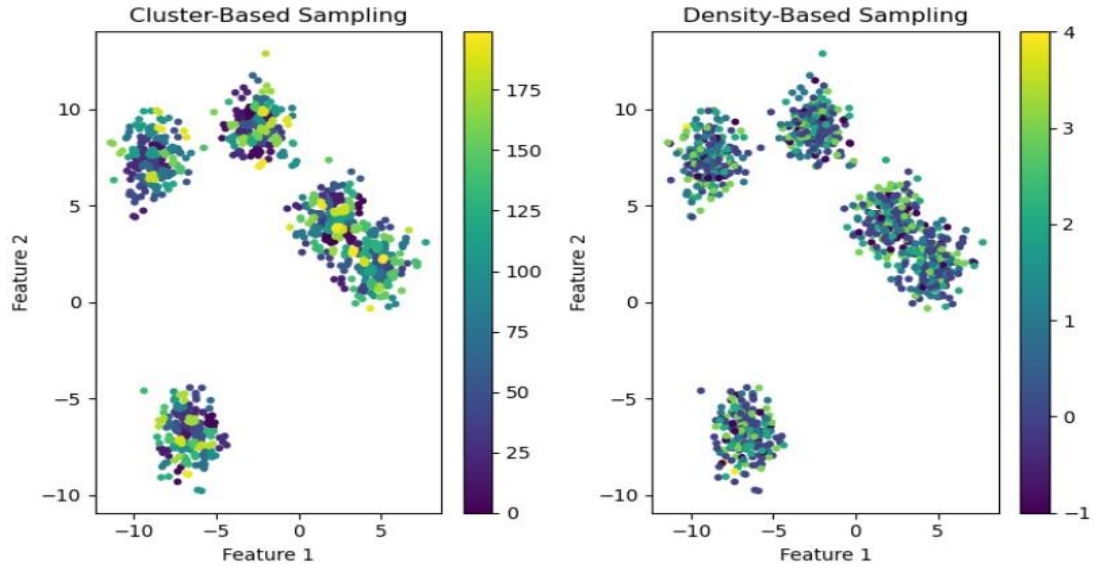
**Figure 5.3 Graph Showing Clustering Pattern Of Cluster-Based And Density-Based Sampling**
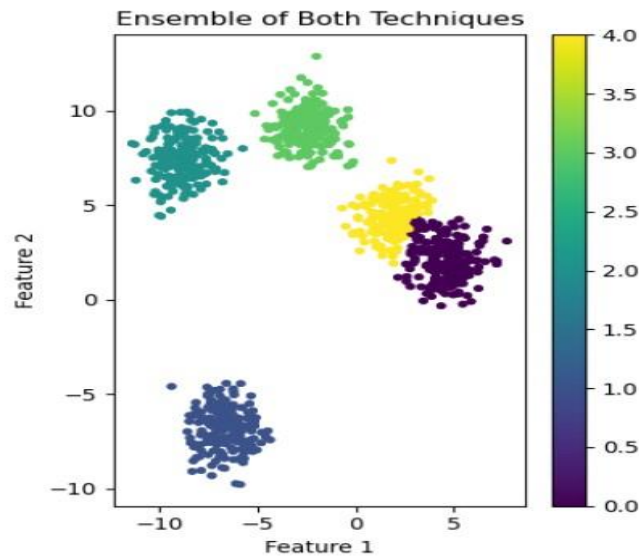


**Figure 5.4 Graph Showing The Clustering Pattern After Ensembling Both Sampling Techniques**

Additionally, the ensemble method also uncovers the new openings for research and innovation in machine learning and data analysis outside of its direct impact on clustering performance as shown in Figure 5.4. This information also highlights the importance of the investigation of the hybrid approaches while developing algorithms, since numerous sampling strategies can be successfully integrated. The performance and adaptability of the ensemble approach can be further enhanced through the utilization of sophisticated

machine learning algorithms or through the examination of other selected sampling procedures within the future research.

In conclusion, the application of our proposed ensemble approach results outlines the comprehensive integration of several sampling approaches to address scalability and quality issues of spectral clustering. Our technique leverages the strengths of both density- and cluster-based sampling, acting as a suitable future path to face the arising demands of big data analysis.

# CHAPTER 6
# CONCLUSIONS AND FUTURE SCOPE

To sum up, we conclude that the ensemble method holds tremendous promise for operating as a viable tool for high-quality clustering of big and complex datasets in a scalable fashion. Our approach unlocks new opportunities in terms of data analysis and interpretation in various fields thanks to the synergy of diverse sampling strategies.

At the same time, several directions imply potential in further exploration and development are visible. First and foremost, the ensemble method can be fine-tuned for even higher performance by implementing various combinations and weights of existing sampling methods. The approach scalability can be improved by leveraging the recent advancements in distributed and parallel computing, which would enable the analysis of even more extensive datasets.

Importantly, our ensemble method is not constrained to traditional clustering applications. It can be modified to operate as recommendation systems, anomaly detection algorithms, and in numerous other pattern recognition tasks. Establishing cooperation with industry and domain experts can bring significant advantages in terms of practical applicability of the proposed approach.

Overall, we conclude that our ensemble method is a significant development in terms of clustering algorithms, offering a flexible and scalable way to address complex datasets. With further research and development.

The ensemble method has the potential to stimulate innovation and simplify data-driven decision-making. However, the practical implications of the ensemble approach are no less vital than its academic achievements. Its capacity to handle vast and complex data quickly and efficiently has a meaningful effect across several essential industries.

Including marketing, finance, healthcare, and cybersecurity. For example, the ensemble approach can help the healthcare sector by discovering patterns in patient information to create individualized treatment protocols and understand approaches of managing

diseases. In the financial industry, detecting unusual activity and clustering related financial transactions can help with fraud detection, creating risk profiles, and optimizing financial portfolios.

In addition, the ensemble approach can significantly impact research in areas such as genomics, where enormous analysis of genomic data is necessary for comprehending complex biological processes and identifying genetic variables causing diseases. The ensemble approach will allow researchers to quickly make progress on personalized healthcare and customized treatment by aiding in the scalable analysis of genetic data. Moreover, the interdisciplinary nature of the ensemble method paves the way for cooperation among data scientists, industry participants, and subject experts. By combining knowledge from many fields, researchers can draw upon the full potential of the ensemble method to tackle complex problems and inspire innovation in data-driven decision-making.

To conclude, the ensemble approach marks a substantial technical shift in clustering algorithms and a harbinger of transformative modifications in numerous fields and industries. Continuing to explore its potential and refine its application will lead to a paradigm shift in the assessment and conclusions drawn from vast and complex datasets, enabling more precise and up-to-date decision-making.

# References

[1] C. Alzate and J. A. K. Suykens, "Hierarchical Kernel Spectral Clustering," IEEE Trans. Pattern Anal. Mach. Intell., vol. 33, no. 8, pp. 1521-1532, Oct. 2011, IEEE.

[2] E. Arias-Castro, G. Chen, and G. Lerman, "Spectral clustering based on local linear approximations," Electron. J. Statist., vol. 5, pp. 1558-1583, 2011, Electronic Journal of Statistics.

[3] R. Couillet and F. Benaych-Georges, "Kernel spectral clustering of large dimensional data," Electron. J. Statist., vol. 11, no. 2, pp. 4011-4049, 2017, Electronic Journal of Statistics.

[4] J. E. Campos, C. de Jesus, H. Ombao, and J. Ortega, "The Hierarchical Spectral Merger algorithm: A New Time Series Clustering Procedure," J. Classification, vol. 33, no. 1, pp. 139-166, 2016, Springer Nature.

[5] W.-Y. Chen et al., "Parallel Spectral Clustering in Distributed Systems," IEEE Trans. Knowl. Data Eng., vol. 23, no. 3, pp. 419-431, Mar. 2011, IEEE.

[6] C. Boutsidis et al., "Randomized Dimensionality Reduction for k-Means Clustering," IEEE Trans. Inf. Theory, vol. 60, no. 9, pp. 4710-4724, Sep. 2014, IEEE.

[7] R. J. Sánchez-García et al., "Hierarchical Spectral Clustering of Power Grids," IEEE Trans. Power Syst., vol. 29, no. 5, pp. 2501-2511, Sep. 2014, IEEE.

[8] H. Jia et al., "The latest research progress on spectral clustering," Neural Comput. Appl., vol. 23, no. 7-8, pp. 2039-2050, Nov. 2013, Springer.

[9] A. Chakeri, H. Farhidzadeh, and L. O. Hall, "Spectral Sparsification in Spectral Clustering," in 2016 23rd International Conference on Pattern Recognition (ICPR), Cancún Center, Cancún, México, Dec. 4-8, 2016, IEEE Xplore Digital Library.

[10] P. Kolev and K. Mehlhorn, "Approximate Spectral Clustering: Efficiency and Guarantees," in 24th Annual European Symposium on Algorithms, Saarland Informatics Campus, Germany, Jul. 29, 2018, Max Planck Institute for Informatics.

[11] H. Sun and L. Zanetti, "Distributed Graph Clustering and Sparsification," in 25th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA), 2018, ACM-SIAM Symposium on Discrete Algorithms.

[12] F. Gao et al., "Distributed Approximate Spectral Clustering for Large-Scale Datasets," in HPDC'12 (ACM Symposium on High Performance Distributed Computing), 2012, ACM.

[13] S. Yoo, H. Huang, and S. P. Kasiviswanathan, "Streaming Spectral Clustering," in IEEE 32nd International Conference on Data Engineering (ICDE), 2016, IEEE.

[14] N. Cristianini, J. Shawe-Taylor, and J. Kandola, "Spectral Kernel Methods for Clustering," in NeurIPS (Conference on Neural Information Processing Systems), 2001, NeurIPS.

[15] D. Yan, L. Huang, and M. I. Jordan, "Fast Approximate Spectral Clustering," Tech. Rep. No. UCB/EECS-2009-45, Department of Statistics, Intel Research Lab, Departments of EECS and Statistics, University of California, Berkeley, 2009, University of California, Berkeley.

[16] N. Tremblay and A. Loukas, "Approximating Spectral Clustering via Sampling: a Review," arXiv:1901.10726 [cs, stat], 2019, arXiv (Cornell University).

[17] D. Achlioptas, "Database-friendly random projections: Johnson-Lindenstrauss with binary coins," Journal of Computer and System Sciences, vol. 66, no. 4, pp. 671–687, 2003.

[18] H. T. Ali and R. Couillet, "Improved spectral community detection in large heterogeneous networks," p. 49.

[19] J. Altschuler et al., "Greedy column subset selection: New bounds and distributed algorithms," in Proceedings of the 33rd International Conference on Machine Learning, M. F. Balcan and K. Q. Weinberger, eds., vol. 48 of Proceedings of Machine Learning Research, New York, New York, USA, June 2016, pp. 2539–2548, PMLR.

[20] E. Anagnostopoulos et al., "Low-quality dimension reduction and high-dimensional approximate nearest neighbor," in 31st International Symposium on Computational Geometry (SoCG 2015), L. Arge and J. Pach, eds., vol. 34 of Leibniz International Proceedings in Informatics (LIPIcs), Dagstuhl, Germany, 2015, pp. 436–450, Schloss Dagstuhl–Leibniz-Zentrum für Informatik.

[21] A. Andoni and P. Indyk, "Near-optimal hashing algorithms for approximate nearest neighbor in high dimensions," in 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), Oct. 2006.

[22] D. Arthur and S. Vassilvitskii, "K-means++: The advantages of careful seeding," in Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, 2007, pp. 1027–1035, Society for Industrial and Applied Mathematics.

[23] S. Arya et al., "An optimal algorithm for approximate nearest neighbor searching fixed dimensions," J. ACM, vol. 45, no. 6, pp. 891–923, Nov. 1998.

[24] M. Aumüller et al., "ANN-benchmarks: A benchmarking tool for approximate nearest neighbor algorithms," in Similarity Search and Applications, C. Beecks et al., eds., 2017, pp. 34–49, Springer International Publishing.

[25] Y. Avrithis et al., "High-dimensional approximate nearest neighbor: k-d generalized randomized forests," arXiv:1603.09596 [cs], Mar. 2016, arXiv:1603.09596.

[26] O. Bachem et al., "Fast and provably good seedings for k-means," in Advances in Neural Information Processing Systems 29, D. D. Lee et al., eds., 2016, pp. 55–63, Curran Associates, Inc.

[27] O. Bachem et al., "Practical coreset constructions for machine learning," arXiv:1703.06476 [stat], Mar. 2017, arXiv:1703.06476.

[28] B. Bahmani et al., "Scalable k-means++," Proc. VLDB Endow., vol. 5, no. 7, pp. 622–633, Mar. 2012.

[29] Z. Bai et al., "Templates for the solution of algebraic eigenvalue problems: A practical guide," SIAM, 2000.

[30] P. C. Bellec et al., "A sharp oracle inequality for graph-slope," Electron. J. Statist., vol. 11, no. 2, pp. 4851–4870, 2017.

[31] Y. Bengio et al., "Label propagation and quadratic criterion," in Semi-Supervised Learning, 2006, pp. 193–216, MIT Press.

[32] R. Bhatia, "Matrix analysis," vol. 169, Springer Science & Business Media, 1997.

[33] T. D. Bie and N. Cristianini, "Fast SDP relaxations of graph cut clustering, transduction, and other combinatorial problems," Journal of Machine Learning Research, vol. 7, Jul. 2006, pp. 1409–1436.

[34] C. Bordenave et al., "Non-backtracking spectrum of random graphs: Community detection and non-regular Ramanujan graphs," in 2015 IEEE 56th Annual Symposium on Foundations of Computer Science, Oct. 2015, pp. 1347–1357.

[35] D. Bouneffouf and I. Birol, "Sampling with minimum sum of squared similarities for Nystrom-based large scale spectral clustering," in IJCAI, 2015, pp. 2313–2319.

[36] C. Boutsidis et al., "Unsupervised feature selection for the k-means clustering problem," in Advances in Neural Information Processing Systems 22, Y. Bengio et al., eds., 2009, pp. 153–161, Curran Associates, Inc.

[37] C. Boutsidis et al., "Spectral clustering via the power method-probably," in International Conference on Machine Learning (ICML), 2015.