

PROBABILISTIC PATH QUERIES IN PATH NETWORKS: EFFECTIVE AND EFFICIENT CLUSTERING METHODS

S. Valenteena Jafflet^{1*}, N. Priya²

¹*Computer Science Engineering Bharath University Chennai, India, ²Assistant Professor, Computer Science Engineering Bharath University Chennai, India

¹jaffletteena@gmail.com, ²shivapriyamari@gmail.com

*Corresponding Author: -

Email ID - shivapriyamari@gmail.com

Abstract: -

Efficiently processing shortest path (SP) queries over stochastic networks attracted a lot of research attention as such queries are very popular in the emerging real-world applications such as Intelligent Transportation Systems and communication networks whose edge weights can be modeled as a random variable. Some previous works aim at finding the most likely SP (the path with largest probability to be SP), and others search the least-expected-weight path. In all these works, the definitions of the shortest path query are based on simple probabilistic models which can be converted into the multi-objective optimal issues on a weighted graph. Challenging problem, two algorithms, the PEEDR and the CPGS clustering algorithm. Reliable clusters are those which are not likely to be disconnected in the context of different instantiations of the uncertain graph. we provide a generalized reliability measurement from two basic intuitions (purity and size balance) to overcome the challenges from standard reliability criterion, and develop a novel k-means algorithm to solve the uncertain clustering problem.

Keywords— Clustering, probabilistic graphs correlated, algorithm



Distributed under Creative Commons CC BY-NC 4.0 OPEN ACCESS

I. INTRODUCTION

Network data analytics lie in the core of many scientific fields such as social, biological and mobile ad-hoc networks. Typically, such network data are associated with uncertainty. This uncertainty is either due to the data collection process or to machine-learning methods employed at preprocessing. Uncertainty may be also added to data for privacy-preserving reasons. We model such uncertain networks as probabilistic graphs. Every edge in a probabilistic graph is associated with a probability of existence. As an example, consider the probabilistic protein-protein interaction (PPI) networks.

Edges in PPI networks represent interactions between proteins which result from combining error-prone measurements and therefore entail uncertainty. We focus on the problem of partitioning a probabilistic graph into clusters. This is a fundamental problem for probabilistic graphs, just as it is for deterministic graphs. Partitioning a probabilistic graph into clusters has many applications such as finding complexes in protein-protein interaction networks and communities of users in social networks.

A straightforward approach to clustering probabilistic graphs is to heuristically cast the probability of every edge into a weight and apply existing graph-clustering algorithms on this weighted graph. Approach is a problem and there is no meaningful way to perform such a casting, but also there is no easy way to additionally encode normal weights on the edges. In this case, ignoring such correlations will lead to incorrect results.

A new algorithm namely, PEEDR, which is rather efficient for clustering correlated probabilistic graphs, and several pruning methods for this algorithm.

Another algorithm, as CPGS, for clustering correlated probabilistic graphs based on the spectral clustering algorithm, which can produce better cluster.

In this subsection, we present a novel algorithm called Partially Expected Edit Distance Reduction (PEEDR), for clustering a correlated probabilistic graph named G . illustrating the PEEDR algorithm, we first present a definition. The PEEDR algorithm initializes a cluster with one vertex. Then for each vertex which is adjacent to cluster, and it has been removed into cluster if it reduces the expected edit distance from G to the current cluster graph. These steps are applied until we cannot expand the cluster.

We next choose a vertex from the unclustered vertices and repeat the above procedure generating another cluster. The procedure has been repeated until all of the vertices of G are grouped into clusters. Until getting the final cluster graph. The problem in the above clustering procedure is which vertex to choosing in each iteration. Motivating by the observation that the vertices with higher degrees are more likely to be the centers of clusters, the vertices in G has been sorted in order of descending degrees.

Inspired by the concept of maximum clique, we propose to build a Distance-Probability-Threshold Clique (**DPTC**) centered at the singleton cluster. The intuition behind DPTC is that vertices with small distances and high similarities are likely to be grouped into the same clusters. The clustering process of the *PEEDR* algorithm starts from a local graph and establishes the cluster graph gradually. As vertices will never be separated grouping once to a cluster, it is likely a greedy algorithm. *PEEDR* algorithm will not meet the need for high precision. Besides, there exists no prior information about the number of final clusters. In some applications, graph clustering aims to partition vertices into a certain number of clusters. Based on this observation, one straightforward approach to extending the spectral clustering algorithm for correlated probabilistic graphs works as follows:

- 1) We map conditional probabilities into weights between each pair of adjacent vertices.
- 2) We extend the Dijkstra method to find the K nearest neighbors (K -NN) of each vertex. It enumerating part of the possible world graphs and calculates the probability that a vertex is the K -NN of the others when correlations exist among edges.
- 3) We establish according a Laplacian matrix to the K -NN resulting, and computing the e vectors of it according to the power method.
- 4) Represent the vertices by points in a K -space dimensional, cluster these points with a K -means algorithm. We call the straightforward method of Spectral; it has been be used as a benchmark method in our experiments.
WLAN.

II RELATED WORKS

2.1 Distance Constraint Reach Ability Computation in Uncertain Graphs

Driven by the emerging network applications, queries and data mining uncertain graphs has become increased important. An un structure attributes also contain a large amount of information for clustering purpose [3]. Here, we investigate a fundamental problem concerning on uncertaining graphs, which is call the distance-constraint reach ability (DCR) problem: Given two vertices ofs and t, which is the probability that the distance from s to t is less than or equal to a userdefined threshold d in the uncertaining graph. This problem is NP-hard, we focus on efficiently and accurately approximating DCR online.

The main results include two new estimators for the probabilistic reach ability. One of the types is a HorvitzThomson estimator based on the unequal probabilistic sample scheme, and another is a novel recursive sample estimator, effectively combines a deterministic recursive computation procedure with a sampling process to boost the estimation accuracy. Both estimators can produce smaller variance than the direct sampling estimator, which considers trials to be either 1 or 0. And also present methods to make these estimators more computationally efficient Therefore, the total number of clustering is output part [4].

The comprehensive experiment evaluation on both real and synthetic datasets demonstrates the efficiency and accuracy of our new estimators. we investigate a fundamental research problem in uncertain graphs: the distance-constraint

reachability (DCR) query problem. In a deterministic directing graph, the reachable query, whether which one vertex can reach another one, is the basis for a variety of database (XML/RDF) and network applications (e.g., social and biological networks).

The importance of distance-constraint reach ability (DCR) query is multi-fold. First, DCR query can contribute to a wide range of real-world applications, which is ranging from social network analysis to biological networks to ontology. the trust ranking between any two persons can be formulated as a distance-constraint reach ability problem; and in the protein-protein interaction network, DCR query can be applied to compute the function similarity between two proteins and the chance they belong to a common protein complex.

2.2 K Nearest Neighbors in Uncertain Graphs

Complex networks, which as social, biological, and communication networks, which is often entailed uncertainty, and it, can be modeled as probabilistic graphs. Similarly, to the problem of search similarity in standard graphs, of a basic problem for probabilistic graphs is to efficiently answer k-nearest neighbor queries (k-NN), which is the problem of computing the k closest nodes to some specific node.

In this paper we introduce a framework for processing k-NN queries in probabilistic graphs. We are proposing novel distance functions that extend to the wellknown graph concepts, as shortest paths [1]. In order to computation them in probabilistic graphs; we have design algorithms based on sampling.

During k-NN query processing we efficiently prune the search space using novel techniques [7]. Biological networks constitute one of the main applications of probabilistic graphs. Nodes represent genes and proteins, and edges represent interactions among them. Since the interactions are derived through noisy and error-prone lab exercise, every edge is associated with an uncertain value. In protein-protein interaction networks, possible interactions have been established experimentally for our experiments indicate that our distance functions outperform previously used alternatives in identifying true neighbors in real-world biological data.

We also demonstrate that our algorithms scale for graphs with tens of millions of edges. The problems of computing distance functions and processing k-NN queries are fundamental for probabilistic graphs, just as they are for standard graphs. They serve as primitive operators for tasks such as link prediction, clustering, classification, and graph mining. We present a principle of extension of these problems in the presence of uncertainty and we assess the quality of the proposed distance functions using a real probabilistic protein-protein interaction network. a limited number of pairs of proteins. Identifying protein neighbors in such interaction networks is useful for predicting possible co-complex memberships and new interactions [8]. Thus, k-nearest neighbor queries can be used to provide candidate links, whose validity can be further explored through strenuous biological experiments.

2.3 Representing and Querying Correlated Tuples in Probabilistic Databases

Probabilistic databases have received considerable attention recently due to the need for storing uncertain data produced by many real-world applications. Widespread use of probabilistic databases is hampered by two limitations:

- (1) current probabilistic databases make simplistic assumptions about the data (e.g., complete independence among tuples) that make it difficult to use them in applications that naturally produce correlating of data, and
- (2) most probabilistic databases can only answer a restricted subset of the queries that can be expressed using traditional query language. Addressing both these limitations by proposing a framework that can represent not only probable tuples, also but correlations that present among them. Our proposed framework lend itself to the possible world semantics thus preserving the precise query semantics extant in current probabilistic databases.

We develop an efficient strategy for query evaluation over such probabilistic databases by casting the query processing problem as an inference problem in an appropriately constructed probabilistic graph model. Presenting several optimizations specific to probabilistic databases that enable resourceful query evaluation. validation of our approach by presenting an experimental evaluation that illustrates the effectiveness of our techniques at answering various queries using real and synthetic datasets.

Database research has primarily concentrated on how to store and query *exact* data. This has led to the development of techniques that allow the user to express and efficiently process complex queries in a declarative fashion over large data collection. Unfortunate, real-world applications produce large amounts of *uncertainty of data*. Such cases, database need to do more than simply store and retrieve; they have to help the user sift through the uncertainty and find the results *most likely to be* the answer.

Natural dependencies in the data: Many application domains naturally produce correlated data.

Dependencies while query evaluation: Problem of handling dependencies among tuples arises naturally during query evaluation even when one assumes that the base data tuples are independent.

The major focus of this tutorial is not about introducing a wide range of graph systems to our people. Although, we endeavor to offer perspectives from a variety of standpoints on the goals and the means for developing a general-purpose graph system. We highlight specific challenges posed by the graph data, the hardware constraints for architect design,

the different types of application which needs, and the power of different models of programming that support such needs. The high complexity of the data (i.e., many relationships exist among the data) means applications require flexible data accesses.

One distinguishing characteristics of graphs is that graph accesses have no locality: As we explore a graph, we invoke random, instead of sequential data accesses, no matter how the graph is stored. In other words, any non-trivial graph query will have poor performance if the locality issue is not proper address. For example, relational models or key/value stores can be used to manage graph data, they do not provide efficient query support, because random accesses are achieved through join operations.

The experimental results also show that algorithm *Crochet* is efficient and scalable. In marketing and customer relation management, customer segmentation is an important task, which partitions customers into groups according to their market behavior such as their purchase records and their responses to marketing campaigns. Customer segmentation in a single market has been extensively studied. However, it is interesting and informative to explore *cross-market customer segmentation*, which identifies groups of customers who have similar behavior in multiple markets. Customer groups found in cross-market customer segmentation can be more coherent and more reliable. Consider the market behavior of six customers, Ann, Bunny, Cathy, Deborah, Ellen and Frank, in two markets, the financial product market and the consumer product market. In a specific market, the similarity among customers in market behavior can be modeled as a similarity graph.

III. PROPOSED METHOD

Graph clustering aims to divide data into clusters according to their similarity, and the number of algorithms have been proposed for clustering graphs, such as spectral clustering, pKwik Cluster algorithm, and k-path clustering likewise. Any how little research has performed to develop efficient clustering algorithms for the probabilistic graphs. In a particular way, it becomes more challenging to efficiently cluster probabilistic graphs when correlations are considered. Defining the problem of clustering correlated probabilistic graphs. Solving the challenging problem, proposing two algorithms, namely the PEEDR and the CPGS clustering algorithm. Each of the proposed algorithms, we are developing several pruning techniques to further improve the efficiency. Evaluating the effectiveness and efficiency of our algorithms and pruning methods through comprehensive experiments. In the CPGS cluster algorithm, which a model is proposed transforming each vertex in a correlated probabilistic transformed points in the multi-dimensional space are iteratively clustered by the K-means algorithm. Addition, we developing several optimization strategies to speed up the clustering process. We formally define the problem of clustering correlated probabilistic graphs and investigate related properties. We propose a new algorithm, PEEDR, which is efficient for clustering correlated probabilistic graphs, and several pruning methods for this algorithm. We develop algorithm, CPGS, for clustering correlating probabilistic graphs based on the spectral clustering algorithm, which can produce better cluster results, even though it is less efficient than PEEDR.

The effectiveness of PEEDR:

In this set of experiments, thus evaluate the effectiveness of the *PEEDR* algorithm. The fundamental algorithm may generate different cluster graphs from the other algorithms. In another hand the words, PLB, PTUB and OROF do not affect the effectiveness of the *PEEDR* algorithm. By comparing Random walking with DPTC, we serve the benefits of the random walk method over the baseline Dijkstra method. The vital idea of the Dijkstra method is to enumerate parts of the possible world graphs which calculate the probability that a DPTC is the KNN of another. As the random method of walk avoids the enumeration of the possible world graphs, it in dealing with probabilistic problems. The experimental result illustrates that it performs better than the Dijkstra method on correlated probabilistic graphs. Compare Random and SCDM walk, we can see the benefits of the optimization using the selfcomplementary matrix method The self-complementary matrix method (SCMM) reduces the number of eigenvectors to be calculated compared to directly calculating the *K* eigenvectors, thus improving the efficiency.

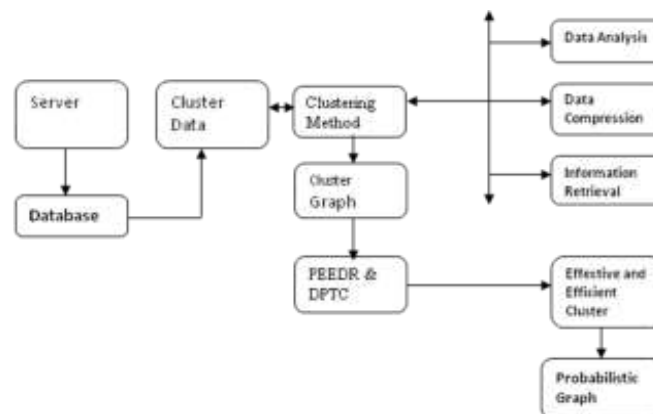


Fig 1. System Architecture

IV. ALGORITHM / METHOD SPECIFICATION

PEEDR Clustering Algorithm:

The PEEDR algorithm initializes a cluster with one vertex. Then for each vertex that is adjacent to cluster, it is removed into the cluster if it reduces the expected edit distance from G to the current cluster graph.

- The step mention is iteratively applied until we cannot expand the cluster.
- We next choose a vertex from the unclustered vertices and repeat the above procedure to generate another cluster.
- The procedure is repeat till all vertices of G are grouped into clusters. We can get the final cluster graph.
- The PEEDR algorithm initializes a cluster with one vertex.
- Then for each vertex that is adjacent to cluster, it is removed the cluster if it reduces to the expected edit distance from node to the current cluster graph.
- We propose to build a Distance-Probability-Threshold Clique (DPTC) centered at the singleton cluster.
- The intuition behind DPTC is that vertices with small distances and high similarities are likely to be grouped into the same clusters.
- The clustering process of the PEEDR algorithm starts from a local graph and establishes the cluster graph gradually. As vertices will never be separated once grouped into a cluster, it is an essential a greedy algorithm.

PEEDR Algorithm:

```

Input: A correlated probabilistic Graph  $G = (V, E, P, F)$ , a
       distance threshold  $d$ , a probability threshold  $\alpha$ 
Output: A Cluster graph Q
1 Sort vertices of G in descending order of their degrees;
2 Initialize  $i \leftarrow 0, b \leftarrow true$ ;
3 Initialize a virtual cluster  $C'$ , where  $V_{C'} \leftarrow V$ ;
4 while ( $V_{C'} \neq \emptyset$ ) do
5     Select the vertex  $v' \in V_{C'}$  with the highest degree;
6     Establish a DPTC (cluster)  $C_i$ , centered with  $v'$  according to  $d$ 
       and  $\alpha$ , and set  $V_{C'} = V_{C'} / V_{C_i}$  // Algorithm 2;
7     while (b = true) do
8          $b \leftarrow false$ ;
9         for ( $v_j \in V_{C'} \cap Adj(C_i)$ ) do
10            If (isReduceEdit( $v_j, C_i$ )) // Algorithm 3 then
11                 $V_{C'} \leftarrow V_{C'} \setminus \{v_j\}, V_{C_i} \leftarrow V_{C_i} \cup \{v_j\}$ ;
12                 $b \leftarrow true$ ;
13            end
14        end
15    end
16     $i \leftarrow i + 1$ ;
17 end
18 return A cluster graph composed of clusters  $C_i (i = 1, 2, \dots, i-1)$ ;
    
```

- CpG cluster uses only arithmetic integer, thus being a fast and computational efficient algorithm able to predict statistically significant clusters of CpG clusters dinucleotides.
- Other of the outstanding feature is that all predicted CGIs start and end with a CpG cluster dinucleotide, which should be proper for a genomic feature whose functionality is based precisely on CpG dinucleotides.
- The parameter of searching in CpGcluster is the distance between two consecutive CpGs cluster, in contrast to prior algorithms.
- Therefore, not any of the main statistical properties of CpG islands (neither G+C content, CpG fraction nor length threshold) are needed as search parameters, which leads to the high specificity and low overlap with spurious Alu elements observed for CpGcluster predictions.

Establishing a DPTC

```

Input: A correlated probabilistic Graph  $G = (V, E, P, F)$ ,
unclustered set  $C$ , an initialized vertex  $x_a$ , a
distance threshold  $d_c$ , a probability threshold  $\alpha$ 
Output: A DPTC which contains vertex  $x_a$ 

1 Initialize a DPTC, where  $V_C \leftarrow \{v_0\}$ ;
2 Initialize  $b \leftarrow true$ ;  $B \leftarrow true$ ;
3 while ( $B = true$ )do
  for ( $v_j \in V_C \cap Adj(C_i)$ ) do
     $B \leftarrow false$ ;  $b \leftarrow true$ ;
4   for ( $v_i \in V_C$ )do
5     if ( $d((x_a, x_i)) \geq D/Q_i$  or  $sim((x_a, x_i)) \leq \alpha$ ) then
6        $b \leftarrow false$ ;
       Break;
     end
   end
   If ( $b=true$ ) then
      $V_C = V_C \cup \{v_j\}$ ;  $V_{C'} = V_{C'} \cup \{v_j\}$ ;
      $B \leftarrow true$ ;
   end
 end
end
Return the DPTC C;

```

Fig.2. Administrator Identification and Password



Fig 3. Update Correlation

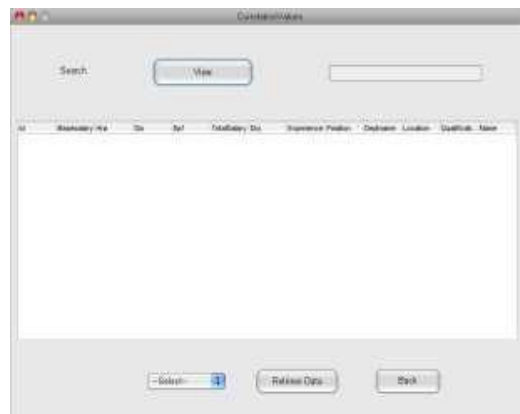


Fig. 5. Pre-Processing

| Id | Username/Pass | Da | Tar | TotalSalary | Exp | Experience | Profile | Department | Location | Qualification | Name |
|-----|---------------|------|-----|-------------|-------|------------|----------|------------|-----------|---------------|----------|
| 190 | 9000 | 508 | 367 | 458 | 0523 | 2013-3-2 | Private | Account | Bangalore | BE | Kumar |
| 191 | 9000 | 389 | 900 | 499 | 7208 | 2012-4-4 | Software | Testing | Chennai | MCA | Senthil |
| 192 | 9000 | 488 | 500 | 444 | 9098 | 2012-6-0 | Software | Architect | Chennai | MCA | Ravi |
| 193 | 9000 | 478 | 345 | 458 | 10281 | 2012-7-0 | Private | Architect | Chennai | BE | Palani |
| 194 | 10000 | 1080 | 230 | 458 | 11889 | 2012-10-4 | Private | Account | Bangalore | BE | Prasanna |
| 195 | 11000 | 1190 | 345 | 453 | 12980 | 2013-11-5 | Private | Account | Bangalore | BE | Ganesh |
| 196 | 12000 | 668 | 456 | 587 | 13829 | 2012-8-5 | Private | Testing | Chennai | MCA | Sharmila |
| 197 | 14000 | 758 | 567 | 782 | 15944 | 2012-2-5 | Teach.L. | Account | Chennai | BE | Ravi |
| 198 | 15000 | 758 | 345 | 458 | 16651 | 2012-6-2 | Teach.L. | Testing | Chennai | MCA | Tarun |
| 199 | 5000 | 289 | 400 | 388 | 9074 | 2010-4-0 | Software | Testing | Chennai | BE | Praga |
| 200 | 9500 | 208 | 400 | 758 | 15480 | 2012-4-6 | Service | Software | Hyderabad | BE | Prasanna |
| 201 | 9500 | 208 | 400 | 758 | 15480 | 2012-4-6 | Service | Software | Hyderabad | BE | Prasanna |
| 202 | 9000 | 458 | 250 | 389 | 10089 | 1989-4-6 | Software | Program | Hyderabad | B.Tech | Srinivas |
| 203 | 9000 | 458 | 300 | 389 | 9850 | 1989-4-4 | Private | Program | Hyderabad | MCA | Raj |
| 204 | 11000 | 458 | 300 | 458 | 12259 | 1989-4-4 | Private | Program | Mumbai | MBA | Dinesh |
| 210 | 9500 | 488 | 450 | 360 | 9778 | 1992-6-0 | Software | Program | Bangalore | BE | Prabhu |
| 216 | 10500 | 488 | 300 | 389 | 11889 | 1990-7-3 | Private | Program | Mumbai | MBA | Krishnak |
| 230 | 12000 | 488 | 450 | 368 | 13990 | 1989-7-4 | Service | Testing | Bangalore | MCA | Ashwin |
| 210 | 9000 | 458 | 300 | 389 | 9850 | 1989-4-4 | Software | Program | Mumbai | MCA | Ravi |
| 210 | 9000 | 458 | 300 | 389 | 9850 | 1989-4-4 | Software | Program | Mumbai | MCA | Ravi |

Fig 6. Correlation Values

V. CONCLUSION

Thus, focused on the problem of finding the cluster graph that minimizes the expected edit distance from the input probabilistic graph. The formulation adheres to the possibleworlds semantics. Also, our objective function does not require the number of clusters as input; the optimal number of clusters is determined algorithmically. We showed that our problem can be efficiently approximated, by establishing a connection with correlation clustering. In addition, we proposed various intuitive heuristics to address it. Further, we established a framework to compute deviations of a random world to the proposed clustering and to test the significance of the resulting clusterings to randomized ones. Also, we addressed versions of our problem where the output clustering is itself noisy. Our experimental evaluation demonstrated that our algorithms not only produce meaningful clusterings with respect to established ground truth, but they also discover the correct number of clusters. Finally, we demonstrated that they scale to graphs with billions of nodes and that they produce statistically significant results. Based on the properties of joint probability, we introduce several pruning methods. A comprehensive performance evaluation verifies the efficiency and effectiveness of our algorithms and pruning methods. This work can be extended to get even better pruning methods and effective algorithm.

REFERENCES

[1].Yu Gu, Member, Chunpeng Gao, Gao Cong, and Ge Yu, VOL. 26, Effective and Efficient Clustering Methods for Correlated Probabilistic GraphsNO. 5, May 2014.

[2].R. Jin, L. Liu, B. Ding, and H. Wang, “Distance-constraint reachability computation in uncertain graphs,” *PVLDB*, vol. 4, no. 9, pp. 551–562, Jun. 2011.

[3].C. C. Aggarwal and H. Wang, *Managing and Mining Graph Data*, New York, NY, USA: Springer, 2010.

[4].M. Potamias, F. Bonchi, A. Gionis, and G. Kollios, “K-nearest neighbors in uncertain graphs,” *PVLDB*, vol. 3, no. 1, pp. 997–1008, Sept. 2010.

[5].X. Lian and L. Chen, “A generic framework for handling uncertain data with local correlations,” *PVLDB*, vol. 4, no. 1, pp. 12–21, 2010.

[6].P. Sen and A. Deshpande, “Representing and querying correlated tuples in probabilistic databases,” in *Proc. ICDE*, Istanbul, Turkey, 2007, pp. 596–605.

[7].A. K. Jain, M. N. Murty, and P. J. Flynn, “Data clustering: A review,”*ACM Comput. Surv.*, vol. 31, no. 3, pp. 264–323, Sept. 1999.

[8].G. Kollios, M. Potamias, and E. Terzi, “Clustering large probabilistic graphs,”*IEEE Trans. Knowl. Data Eng.*, vol. 25, no. 2, pp. 325–336, Feb. 2013.R. Kannan,S. Vempala, and A. Vetta, “On clusterings: Good, bad and spectral,” *J. ACM*, vol. 51, no. 3, pp. 497–515, 2004.



S. Valenteena Jafflet completed M.C.A and her current interested area include data mining and cluster Networking.



N.Priya Assistant Professor in Bharath Universityand her current interested area includes Networking and Cloud Computing.