# Data mining in academic libraries using the K-means algorithm

**Wafaa Razzaq Hashim[1]**          **Serri Ismael Hamad[2]**

[1]University of Thi-Qar – College of Nursing, Iraq
[2]University of Thi-Qar – College of Education for Pure Sciences, Iraq

**Abstract:**

Data mining is an complex process that extracts the required, useful, and comprehensive data from the large quantity of data in order to achieve predetermined goals; some consider data mining to be a common term in the field of extracting knowledge, while others consider data mining to be one of the first steps in the process of extracting knowledge.  And the clustering is the method of separating a data set into distinct cluster groups, each with comparable characteristics, with the goal of combining the data into a single cluster. We are attempting to discuss the possibilities provided by the data mining procedure, and well as how it might improve the standard of service in universities libraries, in this study. Where this article goals to set the books information that is contained in any library, for example, a library College of Education for Pure Science in the University of Thi Qar, by using the K-means clustering approach. Name, Title Of the book, and Author are used as input variables in this K-means clustering method. The result is divided into three groups: (B1) the most commonly borrowed book, (B2) often borrowed book, and (B3) rarely borrowed book. The final result achieved using this K-means Clustering approach includes members of cluster 1 containing (19) members, cluster 2 containing (22) members, and cluster 3 containing (19) members. The library can use the information from grouping this books data. In the case of selecting books to be introduced to the library, it is important to reduce the number of books that are seldom borrowed for the purpose of avoid a backlog of books that are seldom borrowed, leaving room for new books to also be added.

**Keywords:** Data Mining, Cluster, K- means algorithm.

---

## 1. Introduction:

Data mining, also referred as knowledge detection in databases, is the process of examining large data sets in order to find hidden, but potentially useful information. Data mining can uncover hidden

connections and reveal previously undisclosed information trends and patterns. Data mining processes, or models, could be classified according to the mission at hand: association, clustering, classification, and regression [1]. By carefully examining a vast amount of data, data mining reveals the correlations, trends, and patterns of relevant information. Databases, the Web, data warehouses, other data warehouses, or data that is dynamically broadcasted into the system are examples of data sources. By noticing large amounts of data over a period of time, we could indeed deduce previously unknown and helpful data about patterns, designs, trends, and rules in the application domain [2].

In data science, clustering is a beneficial tool. It's a way to determine cluster structure inside a data set based on the highest similarity inside a same cluster and the highest dissimilarity between clusters different. Cluster analysis had become a branch of statistical multi variable analysis, while hierarchical clustering was the original clustering method utilized by biologists and sociologists [3].
Clustering is a relatively common yet crucial subject of research for computer scientists, mathematicians, and pattern matching experts like [4].

Cluster analysis is based on different types of differences of things and uses distance tasks regulations to make the categorization of the sample. If the contributions of points of vectors are clustered, and specimen points in the selfsame collective are centered and specimen points in varied collective are distant, distance functions can be used to classify the points, making statistics in same collective as similar as possible while statistics in various groups are as diverse as possible. Pattern similarity could be measured using the distance function among dots. The metric can be used to identify patterns based on how close the dots are to one another [5, 6].  There are a number of proposed aggregation methods, but K-means is one of the earliest and most often used. The K-means algorithm (Macqueen, 1967) is one of the most essential unsupervised learning algorithms for solution the favor clustering problem. The approach uses an unpretentious and simple approach for specific rating set of data using a predefined number of clusters (assumption k clusters). The central concept is to make k centroids, per cluster. Because different places produce diverse results, these centroids must be positioned carefully. As a result, the ideal choice is to arrange between them as much as possible. The following stride is to take every dot belonging to a specific data set and connect it to the closest centroid. The initial stage is completed, when no points are outstanding. At this stage, we must recalculate k novel centroids to serve as barycenter for the clusters formed in the preceding step. After we've created these k new centroids, we'll need to rebind the conformable data group points to the nearest new centroid. There has been formed a loop, we may note that the k centroids switch their position step by step till no in addition to modifications are made as a product of this episode. In other terms, centroids are no tall moving [4, 7, and 8].

University and college libraries are important places at learning and teaching. The method of servicing these libraries started to change gradually for the development of ideas and knowledge; higher needs are being proposed to work data processing as the quantity of information within libraries continues to rises. Given that user management and provision of services are two of the most important foundations of library management in colleges and universities, the increasing and various demands for services of libraries in the era of large data makes the issues linked to library management and service accuracy increasingly prominent. It is critical that libraries respond quickly to the demands and needs of their patrons and strive to deliver high-quality services [9]. As a result, libraries' first duty should be to identify readers, undertake a reasonable test of their actual demands, and intelligently work on problems and matters connected to their needs. It's also critical for libraries to efficiently manage their readers. Data mining technology can help libraries manage reader services in a more modern and convenient way. Various sorts of data mining are being used in college and university libraries in the age of large data. Without a doubt, the most used clustering approach is the K-mean algorithm. The method was first published by researchers decades ago, and it has since undergone numerous refinements. By minimizing the distance between notes, this algorithm works to determine clusters. The K-means method is used to analyze reading trends in college and university libraries in order to improve service quality and libraries job. It is possible to assist library managers in studying the characteristics of readers and their connections and grouping them into similar sets using the K-means algorithm, which is the first stage in what is known as clustering, in order to know the specifications and real needs to readers. Then, depending on the characteristics and real demands of readers, libraries should design appropriate ways to enhance their services and reader happiness. The k-means algorithm was employed as the clustering technique in this study. The k-means approach will be used to locate the books in a library the University of Thi_Qar's College of Education for Pure Science in order to find the most frequently borrowed books, books that are often borrowed, and books that are seldom borrowed.

## 2. Methodology:

The process of collecting the in most cases borrowed books so at library was done in phases, carrying out the system's analysis, which includes the creation of system input data, such as book lending data from unproven libraries that may be utilized as specimens for analytics. The step of follow is to analyze or handle data entry via clustering methods through the K-means algorithm, and the predictable produced being three clusters: (B1) the more commonly borrowed books, (B2) often borrowed books, and (B3) books that are only seldom borrowed. The data source to be used in the data clustering of this book is borrowed books, which have numerous properties that are needed as a standard for executing the assembly process, as shown in table 1. These attributes are Book Title, Name, and Author.

**Table1.** Data Sampling

| No. | Name | Title of Book | Author | Number of Books |
|---|---|---|---|---|
| 1 | P1 | Web Application Security | Andrew Hoffman | 8 |
| 2 | P2 | Web Application Security | Andrew Hoffman | 8 |
| 3 | P3 | Autocad 14: Instant Reference | Omura, George, Callori, B. Robert | 2 |
| 4 | P4 | Digital Image. Processing | Rafael C. Gonzalez, Richard E. Woods | 2 |
| 5 | P5 | Digital Principles and Applications | Donald P. Leach | 2 |
| 6 | P6 | Digital Principles and Applications | Donald P. Leach | 2 |
| 7 | P7 | Computer Concepts | June Jamrich | 9 |
| 8 | P8 | Computer Concepts | June Jamrich | 9 |
| 9 | P9 | C++ Programming | Malik, D. S. | 19 |
| 10 | P10 | The Pragmatic Programmer | David Thomas. Andrew Hunt | 20 |
| 11 | P11 | clean code | Robert C. | 13 |
| 12 | P12 | Cracking the Coding Interview | Gayle Laakmann McDowell | 19 |
| 13 | P13 | Cracking the Coding Interview | Gayle Laakmann McDowell | 19 |
| 14 | P14 | Cracking the Coding Interview | Gayle Laakmann McDowell | 19 |
| 15 | P15 | Digital Fundamentals | Floyd, Thomas | 15 |
| 16 | P16 | Engineering Electromagnetics | John A. Buck | 5 |
| 17 | P17 | Charles Babbage: Pioneer of the Computer | Hyman, Anthony | 10 |
| 18 | P18 | Microsoft Office Access 2003: The Complete Reference | Anderson, Virginia | 20 |
| 19 | P19 | The Practice of Programming | Brian W Kernighan | 13 |
| 20 | P20 | The Practice of Programming | Brian W Kernighan | 13 |
| 21 | P21 | Genetics and Molecular Biology | Robert Schleif | 2 |
| 22 | P22 | Genetics and Molecular Biology | Robert Schleif | 2 |
| 23 | P23 | Introduction to Molecular Medicine and Gene Therapy | Thomas F. Kresina. Ph.D. | 2 |
| 24 | P24 | Cerebral Signal Transduction | Maarten E. A. Reith | 2 |
| 25 | P25 | Cell Membrane | Yoshihito Yawata | 8 |
| 26 | P26 | Gene Mapping, Discovery, and Expression | Minou Bina | 10 |
| 27 | P27 | Darwin in the Genome | Lynn Helena Caporale | 20 |

# Journal of Education for Pure Science- University of Thi-Qar
## Vol.12, No.1 (March, 2022)
Website: *jceps.utq.edu.iq*                    Email: *jceps@eps.utq.edu.iq*

| 28 | P28 | Conceptual Issues in Evolutionary Biology | Elliott Sober | 13 |
|----|-----|-------------------------------------------|---------------|----|
| 29 | P29 | Applied Cell And Molecular Biology For Engineers | Nindl Waite, Lee Waite | 20 |
| 30 | P30 | Applied Cell And Molecular Biology For Engineers | Nindl Waite, Lee Waite | 20 |
| 31 | P31 | Applied Cell And Molecular Biology For Engineers | Nindl Waite, Lee Waite | 20 |
| 32 | P32 | How to Lie with Statistics | Darrell Huff | 19 |
| 33 | P33 | What is Mathematics? | Richard Courant | 8 |
| 34 | P34 | The Joy of X: A Guided Tour of Math, from One to Infinity | Steven Strogatz | 15 |
| 35 | P35 | The Joy of X: A Guided Tour of Math, from One to Infinity | Steven Strogatz | 15 |
| 36 | P36 | The Joy of X: A Guided Tour of Math, from One to Infinity | Steven Strogatz | 15 |
| 37 | P37 | The Math of Life and Death | Kit Yates | 13 |
| 38 | P38 | Calculus For Dummies | Ryan, Mark | 4 |
| 39 | P39 | Calculus For Dummies | Ryan, Mark | 4 |
| 40 | P40 | Weapons of Math Destruction | Cathy O'Neil | 2 |
| 41 | P41 | Weapons of Math Destruction | Cathy O'Neil | 2 |
| 42 | P42 | The Music of the Primes | Marcus du Sautoy | 8 |
| 43 | P43 | The Music of the Primes | Marcus du Sautoy | 8 |
| 44 | P44 | The Music of the Primes | Marcus du Sautoy | 8 |
| 45 | P45 | Just Enough UNIX | Poul K Andersen | 20 |
| 46 | P46 | Just Enough UNIX | Poul K Andersen | 20 |
| 47 | P47 | Object-oriented and classical software engineering | Stephen R. Schach | 8 |
| 48 | P48 | Control Systems Engineering | Norman Nises | 3 |
| 49 | P49 | Control Systems Engineering | Norman Nises | 3 |
| 50 | P50 | Control Systems Engineering | Norman Nises | 15 |
| 51 | P51 | Every Computer Performance Book | Bob Wescott | 8 |
| 52 | P52 | Plant Biotechnology | Dr. Anupama Tiku, Department of Botany, Ramjas College | 8 |

**Journal of Education for Pure Science- University of Thi-Qar**
**Vol.12, No.1 (March, 2022)**
Website: *jceps.utq.edu.iq*                                        Email: *jceps@eps.utq.edu.iq*

| 53 | P53 | Infinite Powers | Steven Strogatz | 10 |
| 54 | P54 | Infinite Powers | Steven Strogatz | 10 |
| 55 | P55 | The Art of Statistics: How to Learn from Data | David Spiegelhalter's | 5 |
| 56 | P56 | The Art of Statistics: How to Learn from Data | David Spiegelhalter's | 5 |
| 57 | P57 | e: The Story of a Number | Eli Maor | 8 |
| 58 | P58 | e: The Story of a Number | Eli Maor | 8 |
| 59 | P59 | Science and Hypothesis | Henri Poincaré | 5 |
| 60 | P60 | Science and Hypothesis | Henri Poincaré | 5 |

Following the data input has been obtained, the following stage is the data transformed. The data change procedure is carried out at such a step, with the goal of allowing the data to be processed by using the K-means algorithm. Only numbers can be used to run this algorithm. Because the data isn't in a digital format, data elements borrowing books that include titles, names, and authors must be transformed to form of digital. The following are the steps involved in data transformation: Isolate the data by the repeat of its incident; data initialization begins from the uppermost data with a value of 1, then the following data 2, 3, and so on. It can vision data initialization results for variables of authors and title of book in table 2.

**Table2**. Data of Conversion Results

| No. | Title of Book | Author | Number of Books |
|---|---|---|---|
| 1 | 2 | 2 | 8 |
| 2 | 2 | 2 | 8 |
| 3 | 3 | 5 | 2 |
| 4 | 4 | 9 | 2 |
| 5 | 5 | 6 | 2 |
| 6 | 5 | 6 | 2 |
| 7 | 6 | 10 | 9 |
| 8 | 6 | 10 | 9 |
| 9 | 1 | 7 | 19 |
| 10 | 8 | 3 | 20 |
| 11 | 9 | 11 | 13 |
| 12 | 1 | 12 | 19 |
| 13 | 1 | 12 | 19 |
| 14 | 1 | 7 | 19 |
| 15 | 7 | 18 | 15 |
| 16 | 10 | 1 | 5 |
| 17 | 11 | 19 | 10 |
| 18 | 12 | 4 | 20 |
| 19 | 14 | 13 | 13 |
| 20 | 14 | 20 | 13 |

Following the data conversion, the follow stride is to handling data by implementation the k-means clustering algorithm. The first stage is to specify the number of groups that will be generated; in this research the clusters that will be established are the largest number possible three clusters. Then locate of each cluster's initial center point, which is chosen at random. Based on the data, the center dot of each cluster is as next:

The book title for the first center cluster (B1) is 2; the author is 65; the book number is 51; the book title for center cluster (B2) is 25; the author is 16; the book number is 11, and the book title for center cluster (B3) is 40; the author is 35; the number of books is 20.

## 3. Result and Discussion

After obtaining the central dot, the distance between every the data of book and the cluster's center point will be computed till the latter data is the 60th data. Table 3 shows the outcomes of these accounts.

**Table3.** The results of each data's calculation in Iteration 1

| No. | Title of Book | Author | Number of Books | B1 | B2 | B3 | Nearest Distance B1 B2 B3 | | |
|---|---|---|---|---|---|---|---|---|---|
| 1 | | 2 | 8 | 0,0000 | 7,8102 | 19,3391 | 1 0 | 0 | |
| 2 | 2 | 2 | 8 | 0,0000 | 7,8102 | 19,3391 | 1 0 | 0 | |
| 3 | 3 | 5 | 2 | 6,7823 | 2,2361 | 18,0000 | 0 0 | 1 | |
| 4 | 4 | 9 | 2 | 9,4340 | 3,1623 | 14,5945 | 0 0 | 1 | |
| 5 | 5 | 6 | 2 | 7,8102 | 0,0000 | 16,4012 | 0 0 | 1 | |
| 6 | 5 | 6 | 2 | 7,8102 | 0,0000 | 16,4012 | 0 0 | 1 | |
| 7 | 6 | 10 | 9 | 9,0000 | 8,1240 | 10,3441 | 0 0 | 1 | |
| 8 | 6 | 10 | 9 | 9,0000 | 8,1240 | 10,3441 | 0 0 | 1 | |
| 9 | 1 | 7 | 19 | 12,1244 | 17,4929 | 18,0278 | 1 0 | 0 | |
| 10 | 8 | 3 | 20 | 13,4536 | 18,4932 | 19,1050 | 1 0 | 0 | |
| 11 | 9 | 11 | 13 | 12,4499 | 12,7279 | 8,7750 | 0 1 | 0 | |
| 12 | 1 | 12 | 19 | 14,8997 | 18,4662 | 15,1658 | 1 0 | 0 | |
| 13 | 1 | 12 | 19 | 14,8997 | 18,4662 | 15,1658 | 1 0 | 0 | |

| 14 | 1 | 7 | 19 | 12,1244 | 17,4929 | 18,0278 | 1 0 | 0 |
| 15 | 7 | 18 | 15 | 18,1659 | 17,8045 | 6,4807 | 0 1 | 0 |
| 16 | 10 | 1 | 5 | 8,6023 | 7,6811 | 18,7083 | 0 0 | 1 |
| 17 | 11 | 19 | 10 | 19,3391 | 16,4012 | 0,0000 | 0 1 | 0 |
| 18 | 12 | 4 | 20 | 15,7480 | 19,4165 | 18,0555 | 1 0 | 0 |
| 19 | 14 | 13 | 13 | 17,0294 | 15,8430 | 7,3485 | 0 1 | 0 |
| 20 | 14 | 20 | 13 | 22,2036 | 19,9499 | 4,3589 | 0 1 | 0 |

After each of the data is in the nearest cluster, and then recomputed the new cluster center depending on the average organ in the cluster.

The following formula is used to compute the re-generation of a new centroid:

$$B = \frac{\sum m}{d} \tag{3.1}$$

Description: B is the center of data; m is an organ of data that belongs to a specific centroid; and d is the quantity of data that is an organ of a centroid is private. The assembly results are acquired as explained in table 3 based on the minimum value obtained in table 3 above in specifying the centroid value.

**Table4.** Members of every Iteration Cluster 1

| Cluster Name | Cluster Member |
|---|---|
| Cluster 1 | 1, 2, 9, 10, 12, 13, 14, 18, 27, 29, 30, 31, 32, 45, 46, 57, 58 |
| Cluster 2 | 3, 4, 5, 6, 7, 8, 16, 21, 22, 23, 33, 40, 47, 51, 55, 56, 59, 60 |
| Cluster 3 | 11, 15, 17, 19, 20, 24, 25, 26, 28, 34, 35, 36, 37, 38, 39, 41, 42, 43, 44, 48, 49, 50, 52, 53, 54 |

Following the compute above process, the following values will be used to obtain the new centroid. The book title for the initial center cluster (B1) is 5.4706; the books number are 16.8824; the author is 4, 7647; the book title for the center cluster (B2) is 7.7222; the number of books is 4, 6111, the author is 5.5556; and the book title for the center cluster (B3) is 11, 2800; the books number are 9.5200; the author is 19, 3600.

Following the new center dot is acquired, and then does the second repetition with the same computation, viz between the data and the new cluster center. Table 5 shows the computation results of every data to every cluster center dot for the second repetition.

**Table5.** Results of Computation of Every Data to Every Cluster in the 2nd Repetition

| No | Title of Book | Author | Number of Books | B1 | B2 | B3 | Nearest Cluster Distance B1 B3 | B2 |
|----|---------------|--------|-----------------|------|------|------|------|----|
| 1 | 2 | 2 | 8 | 9,9290 | 7,5412 | 19,7433 | 0 0 | 1 |
| 2 | 2 | 2 | 8 | 9,9290 | 7,5412 | 19,7433 | 0 0 | 1 |
| 3 | 3 | 5 | 2 | 15,0879 | 5,4246 | 18,2022 | 0 0 | 1 |
| 4 | 4 | 9 | 2 | 15,5430 | 5,7041 | 14,7268 | 0 0 | 1 |
| 5 | 5 | 6 | 2 | 14,9409 | 3,7981 | 16,5674 | 0 0 | 1 |
| 6 | 5 | 6 | 2 | 14,9409 | 3,7981 | 16,5674 | 0 0 | 1 |
| 7 | 6 | 10 | 9 | 9,4773 | 6,4793 | 10,7591 | 0 0 | 1 |
| 8 | 6 | 10 | 9 | 9,4773 | 6,4793 | 10,7591 | 0 0 | 1 |
| 9 | 1 | 7 | 19 | 5,4284 | 15,9473 | 18,6633 | 1 0 | 0 |
| 10 | 8 | 3 | 20 | 4,3854 | 15,6021 | 19,7038 | 1 0 | 0 |
| 11 | 9 | 11 | 13 | 8,1491 | 10,0821 | 9,3380 | 1 0 | 0 |
| 12 | 1 | 12 | 19 | 8,7647 | 17,1394 | 15,8025 | 1 0 | 0 |
| 13 | 1 | 12 | 19 | 8,7647 | 17,1394 | 15,8025 | 1 0 | 0 |
| 14 | 1 | 7 | 19 | 5,4284 | 15,9473 | 18,6633 | 1 0 | 0 |
| 15 | 7 | 18 | 15 | 13,4557 | 16,2270 | 7,0851 | 0 1 | 0 |
| 16 | 10 | 1 | 5 | 13,2619 | 5,1081 | 18,9515 | 0 0 | 1 |
| 17 | 11 | 19 | 10 | 16,7507 | 14,8505 | 0,6621 | 0 1 | 0 |
| 18 | 12 | 4 | 20 | 7,2758 | 16,0480 | 18,6086 | 1 0 | 0 |
| 19 | 14 | 13 | 13 | 12,4757 | 12,8532 | 7,7433 | 0 1 | 0 |
| 20 | 14 | 20 | 13 | 17,8868 | 17,8445 | 4,4630 | 0 1 | 0 |

**Journal of Education for Pure Science- University of Thi-Qar**
**Vol.12, No.1 (March, 2022)**
Website: *jceps.utq.edu.iq*                                    Email: *jceps@eps.utq.edu.iq*

Computing the distance nearest to the cluster is the next step. Compare members of every repetition cluster 1 and repetition 2, if the cluster organ location still changing, move on to repetition 3. The iteration is terminated if the position of the cluster organ does not change.

Summing values of each cluster individuals and then dividing by the overall cluster number organs are used to specify the center of new cluster dot of the new data.

Based on the aggregation results of each data using the k-means algorithm, the end aggregation results up to the fifth repetition, where the center dot no alteration and no data moves inter clusters. Table 6 shows the cluster results generated of all of these data.

**Table6.** Members of Each Cluster's Results

| Name Cluster | Cluster member |
|---|---|
| Cluster 1 | 9, 10, 11, 12, 13, 14, 18, 27, 28, 29, 30, 31, 32, 34, 35, 36 45, 46, 50 |
| Cluster 2 | 1, 2, 3, 4, 5, 6, 7, 8, 16, 21, 22, 23, 33, 40, 47, 51, 55, 56, 57, 58, 59, 60 |
| Cluster 3 | 15, 17, 19, 20, 24, 25, 26, 37, 38, 39, 41, 42, 43, 44, 48, 49, 52, 53, 54 |

Based on the clusters generated results in table 6, it can be note that the features of cluster 1 individuals are the in most cases borrowed collection of 19 organs, cluster 2 is a collection of books that are predominating borrowed the largest number possible 22 organs, and cluster 3 is a collection of books that is seldom borrowed as many as 19 organs. It can deduce that cluster2 has more organs compared to cluster1 and cluster3 based on the loan data above. So in this status, the library can use this information to select books that should be introduced to the library and to reduce books that are seldom borrowed and there is a spot for books that can be genitive to library.

## 4. Conclusion:

The conclusion that could be extraction from this study is that the K-means clustering algorithm can be used to display the more borrowed books, books recurrent, and books that are very seldom borrowed of the library, consequently that library management can use the information to specify the books to addition to the library and to isolate books that are seldom borrowed for the purpose the places availability. The clustering process results can be used to establish the book site based on the most often borrowed books, which are predominating borrowed, and which are seldom borrowed, allowing students

can readily discovery the required book the campus can also utilize the findings of this clustering operation to a book donation request to the library.

---

**References:**

[1]  Siguenza-Guzman, Lorena, et al. "Literature review of data mining applications in academic libraries." *The Journal of Academic Librarianship* 41.4 (2015): 499-510.

[2]  Huancheng, Liu, Wu Tingting, and Álvaro Rocha. "An analysis of research trends on data mining in Chinese academic libraries." *Journal of Grid Computing* 17.3 (2019): 591-601.

[3]  Sinaga, Kristina P., and Miin-Shen Yang. "Unsupervised K-means clustering algorithm." *IEEE Access* 8 (2020): 80716-80727.

[4]  Haraty, Ramzi A., Mohamad Dimishkieh, and Mehedi Masud. "An enhanced k-means clustering algorithm for pattern discovery in healthcare data." *International Journal of distributed sensor networks* 11.6 (2015): 615740.

[5]  Li, Youguo, and Haiyan Wu. "A clustering method based on K-means algorithm." *Physics Procedia* 25 (2012): 1104-1109.

[6]  Krishna, K., and M. Narasimha Murty. "Genetic K-means algorithm." *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)* 29.3 (1999): 433-439.

[7] Sheshasayee, Ananthi, and P. Sharmila. "Comparative study of fuzzy C means and K means algorithm for requirements clustering." *Indian Journal of Science and Technology* 7.6 (2014): 853.

[8]      Qi, Jianpeng, et al. "An effective and efficient hierarchical K-means clustering algorithm." *International Journal of Distributed Sensor Networks* 13.8 (2017): 1550147717728627.

[9] Tenopir, Carol, Ben Birch, and Suzie Allard. "Academic libraries and research data services: Current practices and plans for the future." (2012).