

SCALING K-MEANS ALGORITHM FOR CLUSTERING LARGE CATEGORICAL DATASETS AND ITS PERFORMANCE ANALYSIS

Rishabh

Arora**

Abstract

Scalable data mining algorithms have become crucial for processing large datasets. The k-means algorithm is one of the most popular algorithms used for clustering. However, its performance degrades significantly when applied to large categorical datasets. In this paper, we propose a modified k-means algorithm which performs much better than the standard k-means algorithm. We conduct an experimental study with various datasets to demonstrate the effectiveness of the proposed algorithm. The experimental results show that the proposed algorithm can reduce the number of iterations required for convergence by up to 50% compared to the standard k-means algorithm, thus improving its scalability.

to efficiently support the clustering of large categorical datasets. The proposed algorithm is designed to handle the high dimensionality and sparsity of categorical data. It uses a novel distance metric to measure the similarity between data points, which is more appropriate for categorical data than the Euclidean distance used in the standard k-means algorithm. The experimental results show that the proposed algorithm achieves better clustering performance than the standard k-means algorithm, especially for large datasets. This makes it a suitable choice for applications requiring efficient clustering of large categorical datasets.

Keywords: Clustering, Categorical Data, K-Means, Scalability

1

Introduction Data mining is the process of extracting useful information from large datasets. It is a multidisciplinary field that combines statistics, computer science, and domain knowledge. One of the primary tasks in data mining is clustering, which involves grouping data points into clusters based on their similarity. The k-means algorithm is a popular clustering method, but it is not well-suited for large categorical datasets. This paper introduces a modified k-means algorithm that addresses the scalability issues of the standard k-means algorithm.

Knowledge discovery in large datasets is a challenging task due to the high volume and complexity of the data. Traditional clustering algorithms like k-means are often inefficient for such datasets. The proposed algorithm aims to overcome these limitations by introducing a more efficient distance metric and a faster convergence mechanism. The experimental results demonstrate that the proposed algorithm is highly scalable and can handle large categorical datasets effectively. This makes it a valuable tool for data analysts working with large-scale categorical data.

1.1 Clustering: Clustering is a fundamental task in data mining that involves partitioning a dataset into groups of similar objects. There are two main types of clustering: hierarchical clustering and partitioning clustering. The k-means algorithm is a partitioning clustering method that iteratively assigns data points to the nearest cluster centroid and updates the centroid until convergence. However, the standard k-means algorithm is not suitable for large categorical datasets due to its high computational complexity.

One of the main challenges in clustering large categorical datasets is the high dimensionality of the data. The proposed algorithm addresses this challenge by using a sparse representation for the data points. This reduces the computational burden and allows the algorithm to scale to larger datasets. Additionally, the proposed algorithm uses a more efficient initialization method to avoid local minima, further improving its performance. The experimental results show that the proposed algorithm is significantly faster and more accurate than the standard k-means algorithm for large categorical datasets.

* SAIT, IIT Bombay, India

** AICTE, IIT Bombay, India

© Technology & Management

© Technology & Management

field of study and relationship among data
biology, economics, and other fields
functional analysis is based on
mathematical models and statistical
analysis such as SPSS and SAS.

about C list analysis
category and in
K-m can K-m order and even the
analysis of a package system

12 Types of C list analysis methods
perform a list of records
Design proposed a list analysis method for

Pairing method These methods construct K -
the list groups with K < <n)

Hierarchical method A hierarchical method creates
given data set This method is hierarchical
perform a list of records (a) Top down (b) Bottom up (c) Both
downward and upward (d) Both downward and upward
method can be used for hierarchical analysis

hierarchical decomposition of
and on how decomposition is
proposed a hierarchical
problem can be addressed

Design based method The general concept is
hierarchical design based on
point

in hierarchical design grows
neighborhood exceeds some limit

Group based method These methods are based on
cluster analysis and advantages of

space partitioning method
list analysis processing

One of the popular list analysis methods is
method hierarchical method
of records with an evaluation and backward K-m
number of records any previous list analysis have
inherent properties and advantages
point of view of hierarchical analysis

is K-m can be based on pairing
basis of evaluation and difference
analysis in order to
focus on number of records because
of difference in number of records
set notation defined.

Design based method is a list analysis method
So a list analysis designed for list analysis use a list
analysis to overcome the problem of category
have used K-m order (a) List analysis with order of category and
have extended B as K-m order of list analysis advanced K -
environment

continuous any category
and order of records
list analysis with order of category
in order of list analysis use a list
analysis with order of category

13 Need of M in list analysis for C list analysis:
possible ways available to find approach
speedup can be achieved by using a list analysis
we can achieve speedup by designing a list analysis
efficient list analysis.

For speedup a list analysis process is
in order of approach. If a
data set is large a list analysis
can be used for list analysis

A list analysis is a list analysis that
using a list analysis to find a list analysis
list analysis and list analysis

list analysis and list analysis
new list analysis point of view
list analysis and list analysis

extent of each point is calculated based on the number of points it is connected to. Thus, the time complexity is $O(n^2)$.

the point is checked by its neighbors. The time complexity is $O(n^2)$.

Proposed Enhancement: In our previous work, we proposed an enhancement to the algorithm. The time complexity is $O(n^2)$.

In our previous work, we proposed an enhancement to the algorithm.

with respect to the time complexity.

only new added points are checked. The time complexity is $O(n^2)$.

Similar to the previous work, we have proposed an enhancement to the algorithm. The time complexity is $O(n^2)$.

the point is checked by its neighbors. The time complexity is $O(n^2)$.

Justification: The main idea of the algorithm is to find all points that are connected to each other. The time complexity is $O(n^2)$.

at $O(n^2)$ time.

distance and minimum distance. The time complexity is $O(n^2)$.

This paper is organized as follows. Section 2 describes the background and problem. Section 3 describes the basic K-M algorithm. Section 4 describes the proposed algorithm. Section 5 presents the experimental environment and results. Section 6 presents the conclusion.

be used to find the minimum distance. The time complexity is $O(n^2)$.

7 present the performance analysis of the proposed algorithm.

Related Work: Many algorithms have been proposed for finding the minimum distance between points. The time complexity is $O(n^2)$.

of the result.

has been used extensively in many areas such as dynamic programming, graph algorithms, and image processing. The time complexity is $O(n^2)$.

Umut A. Cakir et al. [1] have proposed an algorithm for finding the minimum distance between points. The time complexity is $O(n^2)$.

the time complexity is $O(n^2)$.

ionefontinchem em oizomilboup
 condheporm ancofin em oizoml anysubt
 cofoqualcheckingandhecaheepicem ent
 hel fencobtw eexponentilndheannni

eform andheusem usbeabto
 esofm em oizomiludinghe
 polyform em othebanm ake
 ngin e.

G andiposedtheCA CTU Sllhem rhodo
 sum m aisTo fid outhe distrhey peform he
 sum m aindusingandvallaon
 sum m aynform abnform hedatshhedie
 inform abndicoversofandiladiter
 heactakofluseform hesofandilad
 hedeftion ofadistrw henhedatconsis
 induced a ftsum m aizon based aloghm for
 categorallat.

bgyfodistngbyusingulster
 process in hree phases
 hhesum m aizonphasitorm putshe
 ingphasduseshesum m ay
 hhevalaizonphasitorm ies
 edistO verCA CTU Sform aded
 ofcategoricalbutsandhen
 dicoverng such distrin

C lusingofcategoricalanextensivelye
 anddatabaseresearchbutbyohedipl
 categorallatH anedll addresshepro
 m akabasketlatabeypresenngfrequent
 hypegraphThew eghtthegraphicorm puts
 poshtasocionushatanegeneratdf
 pairingalghm iem phytopanhei
 hypegedlthegolhm doesnoproduceit
 obvioushow tobtionofm hem dista
 [14]foracategoricalapproachtoit

eachedaeanoonlybydatm ing
 inepeopleinheaeofdistng
 blm ofdistngansoonina
 m setshypeadgenaw eghtd
 haveageficconfidencd
 om hem sethenahypegraph
 em sn nim inghew eghthecut
 ingofansconsndhot
 nohealtqapebyG lsonat
 basedmonheardynam icayem .

3D efimO ftem SAndBasForm ulon

31C atgorialD ata A refedihapeacategoricalm e
 w hihhaveonlycategoricalbutH ew econs
 categoralbutandmonislem ulal

andlatobjts
 dealngvaluedabutas
 uedbutascategoricalbut

32C atgorialDom ansandA tbutts Let $X_1, X_2, X_3, \dots, X_n$ be hed categorical
 abutsdefingaspae $Q(\theta)$ and $\text{Dom}(X_1), \text{Dom}(X_2), \dots, \text{Dom}(X_n)$ hedom anof
 abut $\text{Dom}(X_1)$ hedatposhtvalshhebu e X_1 andmorleed.

33C atgorialD bjts Let R be acategoricalbjt $R, R_2, R_3, \dots, R_n, \dots, R_\infty$ w e
 canlefe

$R_i = R_j$ (EquivalentC atgorialbjts $overlpm rhod$)

If $R_k = R_j$ for $k = 1, 2, \dots, d$,

Inhiproctw ehavedistecordsbasedon
 overlpm easionelhem . datdiven sin lly m easues

3 Similarity Measure:

similarity measure between two objects of the same category or two objects of different categories.

Let X and Y be two categorical objects. Then X and Y can be defined as follows: $X = \{x_1, x_2, \dots, x_n\}$ and $Y = \{y_1, y_2, \dots, y_m\}$ where x_i and y_j are the values of the attributes.

he
ve

Not that have converted similarity measure into a distance measure. In order to make the similarity measure a distance measure, we have to convert the similarity measure into a distance measure. This can be done by using the following formula:

initial proposed distance
some part in history. The
next of similarity distance

$$sim = \frac{1}{1 + dist}$$

A similarity measure assigns a value between two objects X and Y belonging to the same class or different classes.

value between two objects X

$$S(X, Y) = \sum_{k=1}^d w_k S_k(X_k, Y_k)$$

Where $S_k(X_k, Y_k)$ is the similarity between two objects X_k and Y_k in the k th attribute. w_k is the weight assigned to the k th attribute.

between two values for the
□ A K-Nearest Neighbors algorithm

4 K-Means Algorithm

4.1 Basic K-Means Algorithm:

1. Initialize the number of clusters K and the initial cluster centers $\mu_1, \mu_2, \dots, \mu_K$.
2. Assign each data point to the cluster whose center is closest to it.
3. Calculate the new cluster centers $\mu_1, \mu_2, \dots, \mu_K$ as the mean of all data points assigned to each cluster.
4. Repeat steps 2 and 3 until the cluster centers no longer change significantly.

Sequential K-Means algorithm works by assigning

K clusters. In each iteration, the algorithm assigns each data point to the cluster whose center is closest to it. Then, the new cluster centers are calculated as the mean of all data points assigned to each cluster. This process is repeated until the cluster centers no longer change significantly.

K

Data

Count of similarity distance

P

Number of clusters (k)
Data set (D)

Process

Step 1:

Select initial cluster centers $\mu_1, \mu_2, \dots, \mu_K$ and assign them as initial cluster centers.

and assign them as initial

Step 2:

Select one category label point object
 he proposed this sign the new h
 current point form a/
 For every $R_i \in D(S)$
 {
 For $j = I_k = I_{k-1} + 1$ to I_k
 $Count \leftarrow \text{Sim}(Y_{j-1}, R_i)$
 While Y_j is current object in list
 Index = Indexnum be list to hom
 Add R_i object index and calculate the
 }
Step3 : While object is added to R check
 object is current object in list
 with this new object in the new dis
 tance between category object
Step4 : Repeat step 3 until object has changed list
 data

4.2 Advanced K - m mode Algorithm for Categorical Data
 m mode algorithm uses technique of m - o - n - W
 data point in the cluster

Data

P

Process

Step1 :

Step2 :

Select one category label point object
 he proposed this sign the new h
 current point form a/
 For every $R_i \in D(S)$
 {
 For $j = I_k = I_{k-1} + 1$ to I_k
 $Count \leftarrow \text{Sim}(Y_{j-1}, R_i)$
 While Y_j is current object in list
 Index = Indexnum be list to hom

m be data and assign
 o - m - o - d - e - m - o - s - i - m - i - r - h
 lists
 R_i in o - s - i - m - i - r
 item mode
 item is not categorical
 categorical object found
 calculate the current
 & calculate the mode
 finally select the
Object The advanced K -
 m - o - d - e - m - o - s - i - m - i - r - h

```

    }
    Add  $R_i$  object index list
    Calculate next flat point
    Merge  $R_i$  and  $R_{i-1}$ 
    Recalculate  $R_i$ 
  }

```

Step3 :

```

  Repeat 3 until no change
  For  $k = 1$  to  $n$ 
  {
    For  $i = 1$  to  $n$ 
    {
      For  $j = 1$  to  $n$ 
      {
         $Count \leftarrow \text{Sum}(R_i, M_j)$ 
         $Index \leftarrow \text{Index of min value}$ 
        {
          Remove( $R_i, j$ )
          Add( $R_{Index}$ )
        }
      }
    }
  }

```

Step4 :

```

  Repeat 3 until no change

```

Advanced K-modes algorithm has following properties

1. Basic K-modes algorithm has complexity of $O(km \cdot n^2)$ where k is number of clusters, m is number of categorical attributes and n is number of objects. Advanced K-modes algorithm has reduced complexity by using an efficient new algorithm having complexity of $O(kmn)$.
2. A stable algorithm of any form.
3. Because of its simplicity and efficiency.

To illustrate K-modes algorithm and advanced version, we assume that there are 1000 categorical objects with 22 categorical attributes. Then the complexity of sequential algorithm will be $1000 \times 1000 \times 4 \times 22 \times t \approx 88000000t$ where t is the number of iterations. In the advanced K-modes algorithm, the complexity is $1000 \times 22 \times t \approx 22000t$.

Now we have applied the technique to the K-modes algorithm. The complexity of the advanced K-modes algorithm is $1000 \times 22 \times t \approx 22000t$ where t is the number of iterations. The complexity of the sequential algorithm is $1000 \times 1000 \times 4 \times 22 \times t \approx 88000000t$.

Experiment Environment

5.1.1. Dataset We have used M ushroom dataset which contains putative binary labels for each item. We have evaluated the performance of both the advanced and the baseline algorithms. We have used the paper [1] as a reference.

The M ushroom dataset is a standard dataset for recommendation systems. It has been used by many researchers in the field of recommendation systems. We have chosen this dataset because it is a standard dataset for recommendation systems and it is a standard dataset for recommendation systems.

The M ushroom dataset has 8124 observations and 4208 observations in the poisonous class. Each observation is identified by a unique ID. The dataset is divided into two parts: a training set and a test set. The training set contains 4208 observations and the test set contains 3916 observations.

5.1.2. Evaluation We used the K-fold cross-validation method to evaluate the performance of the algorithms. The evaluation metrics used are the Area Under the Curve (AUC) and the Precision-Recall curve. The evaluation metrics used are the Area Under the Curve (AUC) and the Precision-Recall curve.

5.1.3. Environment The experiment is performed on a Windows 10 machine with 8GB RAM and Intel Core i7 processor. The operating system is Windows 10. The hardware configuration is P4, 512 MB RAM. The software configuration is Java 8, Eclipse IDE.

Conclusion

6.1.1. Changing the number of records We have changed the number of records from 500 to 1000. We have evaluated the performance of the algorithms for different numbers of records. We have evaluated the performance of the algorithms for different numbers of records.

Table 1 shows the performance of the algorithms for different numbers of records. The table shows the performance of the algorithms for different numbers of records. The table shows the performance of the algorithms for different numbers of records.

Table 1: Performance of the algorithms for different numbers of records

N of records	Time to execute Basic implementation (in milliseconds)	Time to execute Advanced implementation (in milliseconds)
500	2574	1748
600	3245	2043
700	5975	2402
800	7332	2761
900	10171	2995
1000	11638	3572

Figure 1: Comparison of response time.

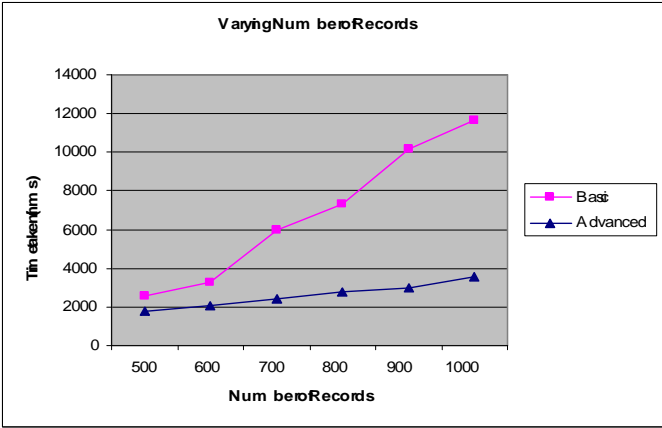


Figure 2: Variation of Time

Change in Number of Clusters: We have fixed the number of records at 1000 and changed the number of clusters from 6 to 10 and evaluated the execution time by both algorithms and draw a graph between number of clusters and execution time (Figure 2).

Table 2 shows the comparison of number of clusters by Basic and Advanced algorithms:

Table 2: Variation of Time with Number of Clusters

N of A butes	Time to execute Basic in permutation (milliseconds)	Time to execute Advanced in permutation (milliseconds)
4	4337	2044
7	4181	2402
10	7160	2917
13	5506	2418
16	6817	2996
19	7129	2699
22	12511	3666

Figure 3: Comparison of performance between basic and advanced algorithms.

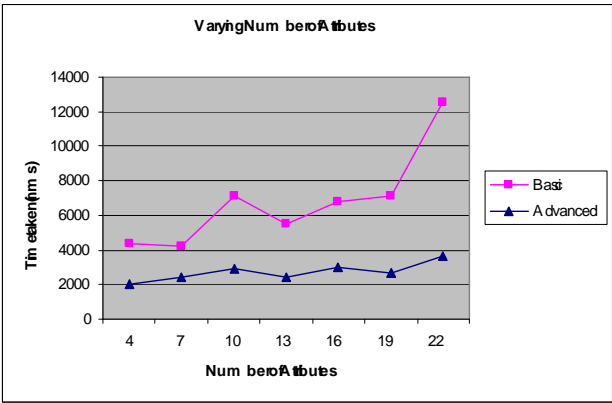


Figure 3: Comparison of performance between basic and advanced algorithms.

We have evaluated the decrease in time in Table 1a considering the number of attributes and the number of attributes in the dataset. We have kept the number of attributes and the number of attributes in the dataset constant and evaluated the time by changing the number of attributes. The results are shown in Table 1a and Figure 3. We observe that the advanced algorithm is faster than the basic algorithm.

Figure 1 shows the corresponding figure with high number of attributes and change in the number of attributes. We have evaluated the time by changing the number of attributes. The results are shown in Table 1a and Figure 3. We observe that the advanced algorithm is faster than the basic algorithm.

Conclusion The Advanced version of K-model algorithm presented in this paper has complexity of $O(mn)$ and achieves better memory organization technique. Such algorithm is capable of performing heap sorting on heterogeneous data set of categorical objects in a number of small and medium size manageable objects. In particular, it shows that the clustering problem can be optimized by adopting memory organization method.

[illegible]

References

- | | | | |
|-----|--|--|------------------------------|
| 1. | U m A c G u y E B a d R o b b H a p S e d i
January 15/2003, New Orleans, Louisiana, USA | veM em oir of POI 03, | |
| 2. | Lukas Zick K C S i a m a k i n n a S u r e s h g a n n a h a
Concurrency and Com m u n i c a t i o n C F P 0 9 A u g u s t 1 - 5
Edinburgh, Scotland, U K . | r, a l l M em oir of
q t m b e 2009, | |
| 3. | Raymond L N g n d i w H a C L A R A N S A M d h o d r
f b p a d a t M i n g E E T a n s a c i o n O n k n o w l
Engineering d 4 N o S e p t e m b e r O c t o b e 2002 | C l i n g O b j e c t
e d g e A n d a t | |
| 4. | Shyam B o i t y a n C h a n d o l a n d V i n k u m a s i n i
C a t e g o r i d a t A C o m p a r e E v a l u a t i o n P r o c
M i n g C o n f e r e n c e A p r i l 2008 a n d G A . | a t M e a s u r e f o r
e e d i n g s 2008 S I A M D a t | |
| 5. | V G a d i c h a k e R a m a k i n n a C A C T U S C i e
U s i n g S u m m a r i s k n o w l e d g e D i c o v e r i d a t M i n g | i n g C a t e g o r i d a t
1999. | |
| 6. | A A h m a d a n d L D a y m d h o d
c a t e g o r i e s f a m e a b o u t
d a t a R e c o g n i t i o n L e t | c o m p u t e d i n d b e t w e e n
i n s u p e r v i s e d t a r i n g f o r
28(1) 110-1 18/2007. | t w o
c a t e g o r i a l |
| 7. | M a i C a m i n B a i t i l u m b e t L R a z e n A g m a
T r i n a r e f i n a p p r o a c h t o c l u p k e m o d e
d a t a b a s e X X S i m p o s i B a n d B a n c o d e d a d o | M T r i n C a t i n o
b a s e d h o l m i n g e
2003 | |
| 8. | M E s t r i k e g l a n d e X X u t A D e n s i
D i c o v e r i n g C l u s t e r L a g e S p a d a t a b a s e P r
K n o w l e d g e D i c o v e r a n d d a t M i n g K D D 96 p p 226
1996. | t y b a s e d h o l m f o r
o d 99 d i n f
23 I P o l n d O R A u g | |
| 9. | S u d i p t G u h a R a j e R a o K y u s o k S h i n C U R E A
a g e l a n d P r o c e e d i n g s I E A C M S i g m o d i t
d a t 1998. | d i s t i n g h o l m f o r
c o n f e r e n c e o n m a n a g e m e n t f
d a t 1998. | |
| 10. | D a v i d G h o s t o n K t h e g P r i b h a k a R a g h v a n i C l u
A n a p p r o a c h B a s e d o n d y n a m i c S y s t e m ' P r o d 998
D a t a b a s e p p 1132 N e w Y o r k A u g u s t 1998. | s i n g C a t e g o r i d a t
d i n f o n d n v a y L a g e | |

11. Y Zhang, A daW a, C haF u, C hunH ing, P engA nn
C ategoriD atI Pro2000IEEEI onD a
U SAM ach2000.
12. ZH uang, ExnsiortheK M eanA golum fid
wIC ategori/ aliesD atM ingandK now edged
1998.
13. CLGB sw asU nsupervied Learningw fM ied
D atEEET nsatornK now edgedandD atEngn
2002.
14. TChiD Fang, ChenY ang, ARouband
A golum fidM iedTypeA butsiL agD atba
2001I onD nK now edged icoveryandD atM i
15. K aufm aL andR ouseauw H 2005Findnggr
introduintulstana, JohnW ileyandSo

H engC lustering
EngineeringSanD eigo,

usingL agD atas
icoverypp283304,

Num ericandN om inal
eering ol4N el,

ScalibC lustering
sEnvironm entProc.
ningpp263268,2001.
oupdata n

ns