

Local Dimensionality Reduction: A New Approach to Indexing High Dimensional Spaces *

Kaushik Chakrabarti
University of Illinois
kaushikc@cs.uiuc.edu

Sharad Mehrotra
University of California
sharad@ics.uci.edu

Abstract

Many emerging application domains require database systems to support efficient access over highly multidimensional datasets. The current state-of-the-art technique to indexing high dimensional data is to first reduce the dimensionality of the data using Principal Component Analysis and then indexing the reduced-dimensionality space using a multidimensional index structure. The above technique, referred to as global dimensionality reduction (GDR), works well when the data set is globally correlated, i.e. most of the variation in the data can be captured by a few dimensions. In practice, datasets are often not globally correlated. In such cases, reducing the data dimensionality using GDR causes significant loss of distance information resulting in a large number of false positives and hence a high query cost. Even when a global correlation does not exist, there may exist subsets of data that are locally correlated. In this paper, we propose a technique called Local Dimensionality Reduction (LDR) that tries to find local correlations in the data and performs dimensionality reduction on the locally correlated clusters of data individually. We develop an index structure that exploits the correlated clusters to efficiently support point, range and k-nearest neighbor queries over high dimensional datasets. Our experiments on synthetic as well as real-life datasets show that our technique (1) reduces the dimensionality of the data with significantly lower loss in distance information compared to GDR and (2) significantly outperforms the GDR, original space indexing and linear scan techniques in terms of the query cost for both synthetic and real-life datasets.

1 Introduction

With an increasing number of new database applications dealing with highly multidimensional datasets, techniques to support efficient query processing over such data sets has emerged as an important research area. These applications include multimedia content-based retrieval, exploratory data analysis/data mining, scientific databases, medical applications and time-series matching. For example, in multimedia retrieval, the objects (e.g., images) are represented by their features (e.g., color histograms, texture vectors and shape descriptors) which define high dimensional feature spaces (HDFS) [17, 37]. In data mining applications, objects are represented by several numeric attributes which again define a HDFS over which the data mining task (e.g., clustering, classification) is performed [3, 33]. HDFSs are also becoming increasingly common in scientific (e.g., SDSS's astronomy database [41]) and medical databases [30]. To provide efficient access over HDFSs, many indexing techniques have been proposed in the literature. One class of techniques comprises of *high dimensional index*

*This work was supported by NSF CAREER award IIS-9734300, and in part by the Army Research Laboratory under Cooperative Agreement No. DAAL01-96-2-0003.

Report Documentation Page

Form Approved
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE 2000		2. REPORT TYPE		3. DATES COVERED 00-00-2000 to 00-00-2000	
4. TITLE AND SUBTITLE Locally Dimensionality Reduction: A New Approach to Indexing High Dimensional Spaces				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) University of Illinois at Urbana-Champaign, Department of Computer Science, 201 N. Goodwin Avenue, Urbana, IL, 61802-2302				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES 26	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

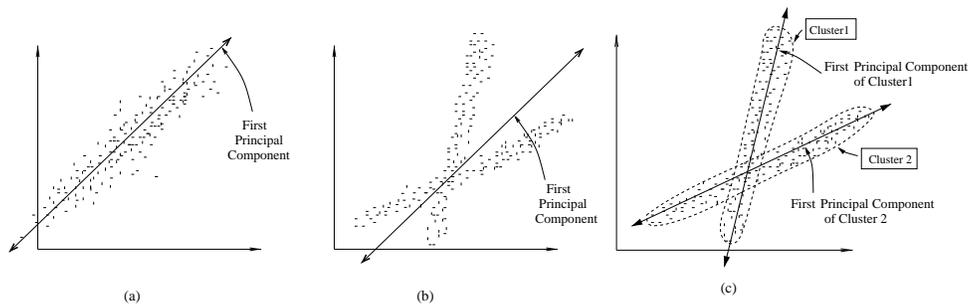


Figure 1: Global and Local Dimensionality Reduction Techniques (a) GDR(from 2-d to 1-d) on globally correlated data (b) GDR (from 2-d to 1-d) on globally non-correlated (but locally correlated) data (c) LDR (from 2-d to 1-d) on the same data as in (b)

trees [5, 44, 28, 11, 31, 7]. Although these index structures work well in low to medium dimensionality spaces (upto 20-30 dimensions), a simple sequential scan usually performs better at higher dimensionalities [6, 43].

To scale to higher dimensionalities, a commonly used approach is *dimensionality reduction* [20]. This technique has been proposed for both multimedia retrieval [17, 36, 27, 42] and data mining ([18, 4, 21]) applications. The idea is to first reduce the dimensionality of the data and then index the reduced space using a multidimensional index structure [17]. Most of the information in the dataset is condensed to a few dimensions (the first few principal components (PCs)) by using principal component analysis (PCA). The PCs can be arbitrarily oriented with respect to the original axes (see Appendix A for details on PCA). The remaining dimensions (i.e. the later components) are eliminated and the index is built on the reduced space. To answer queries, the query is first mapped to the reduced space and then executed on the index structure. Since the distance in the reduced-dimensional space lower bounds the distance in the original space, the query processing algorithm can guarantee no false dismissals [17, 16]. The answer set returned can have false positives (i.e. false admissions) which are eliminated before it is returned to the user. We refer to this technique as *global dimensionality reduction* (GDR) i.e. dimensionality reduction over the *entire* dataset taken together.

GDR works well when the dataset is *globally correlated* i.e. most of the variation in the data can be captured by a few orthonormal dimensions (the first few PCs). Such a case is illustrated in Figure 1(a) where a single dimension (the first PC) captures the variation of data in the 2-d space. In such cases, it is possible to eliminate most of the dimensions (the later PCs) with little or no loss of distance information. However, in practice, the dataset may not be globally correlated (see Figure 1(b)). In such cases, reducing the data dimensionality using GDR will cause a significant loss of distance information. Loss in distance information is manifested by a large number of false positives and is measured by precision [27] (cf. Section 5). More the loss, larger the number of false positives, lower the precision. False positives increase the cost of the query by (1) causing the query to make unnecessary accesses to nodes of the index structure and (2) adding to the post-processing cost of the query, that of checking the objects returned by the index and eliminating the false positives. The cost increases with the increase in the number of false positives. Note that false positives do not affect the quality the answers as they are not returned to the user.

Even when a global correlation does not exist, there may exist subsets of data that are *locally correlated* (e.g., the data in Figure 1(b) is not globally correlated but is locally correlated as shown in Figure 1(c)). Obviously, the correlation structure (the PCs) differ from one subset to another as otherwise they would be globally correlated. We refer to these subsets as *correlated clusters* or simply *clusters*.¹ In such cases, GDR would not be able to

¹Note that correlated clusters (formally defined in Section 3) differ from the usual definition of clusters i.e. a set of spatially close

obtain a single reduced space of desired dimensionality for the entire dataset without significant loss of query accuracy. If we perform dimensionality reduction on each cluster *individually* (assuming we can find the clusters) rather than on the entire dataset, we can obtain a set of different reduced spaces of desired dimensionality (as shown in Figure 1(c)) which together cover the entire dataset² but achieves it with minimal loss of query precision and hence significantly lower query cost. We refer to this approach as local dimensionality reduction (LDR).

Contributions: In this paper, we propose LDR as an approach to high dimensional indexing. Our contributions can be summarized as follows:

- We develop an algorithm to discover correlated clusters in the dataset. Like any clustering problem, the problem, in general, is NP-Hard [32]. Hence, our algorithm is heuristic-based. Our algorithm performs dimensionality reduction of each cluster individually to obtain the reduced space (referred to as subspace) for each cluster. The data items that do not belong to any cluster are outputted as outliers. The algorithm allows the user to control the amount of information loss incurred by dimensionality reduction and hence the query precision/cost.
- We present a technique to index the subspaces individually. We present query processing algorithms for point, range and k-nearest neighbor (k-NN) queries that execute on the index structure. Unlike many previous techniques [27, 42], our algorithms guarantee correctness of the result i.e. returns exactly the same answers as if the query executed on the original space. In other words, the answer set returned to the user has no false positives or false negatives.
- We perform extensive experiments on synthetic as well as real-life datasets to evaluate the effectiveness of LDR as an indexing technique and compare it with other techniques, namely, GDR, index structure on the original HDFS (referred to as the original space indexing (OSI) technique) and linear scan. Our experiments show that (1) LDR can reduce dimensionality with significantly lower loss in query precision as compared to GDR technique. For the same reduced dimensionality, LDR outperforms GDR by almost an order of magnitude in terms of precision. and (2) LDR performs significantly better than other techniques, namely GDR, original space indexing and sequential scan, in terms of query cost for both synthetic and real-life datasets.

Roadmap: The rest of the paper is organized as follows. In Section 2, we provide an overview of related work. In Section 3, we present the algorithm to discover the correlated clusters in the data. Section 4 discusses techniques to index the subspaces and support similarity queries on top of the index structure. In Section 5, we present the performance results. Section 6 offers the final concluding remarks.

2 Related Work

In this section, we discuss the related work on high dimensional index structures, global dimensionality reduction and clustering algorithms.

High Dimensional Index Structures Recent research on high dimensional indexing has led to the development of several index structures including X-tree[5], SS-tree [44], SR-tree [28], M-tree [11], TV-tree [31] and Hybrid-tree [7]. These index structures use novel data/space partitioning strategies and scale better to high dimensionalities compared to spatial index structures (e.g., R-tree, grid file). They are extensively used for similarity search in multimedia retrieval [17, 10], data mining [14, 3] and decision support [40, 13] applications. Although

points. To avoid confusion, we refer to the latter as *spatial clusters* in this paper.

²The set of reduced spaces may not necessarily cover the entire dataset as there may be outliers. We account for outliers in our algorithm.

these index structures can scale to medium dimensionalities (upto 20-30 dimensions), above a certain dimensionality (referred to as the critical dimensionality), they are outperformed by a simple sequential scan through the database [43, 6]. The reason is that the data space becomes sparse at high dimensionalities causing the bounding regions to become large. The query ends up overlapping with most nodes of the tree resulting in a large number of disk accesses and hence a high query cost. The linear scan performs better in such cases since sequential I/O is significantly cheaper compared to random I/O. Obviously, the critical dimensionality depends on the dataset and the index structure used.

Global Dimensionality Reduction GDR techniques has been studied extensively in statistical pattern recognition and multivariate data analysis. The principal component analysis (PCA) or Karhunen-Loeve (K-L) transform is the optimal way of mapping points in a D -dimensional space to points in a d -dimensional space ($d \leq D$) [12, 20]. The mapping is optimal in the sense it minimizes the mean square error (MSE), where the error is the distance between each D -d point and its d -d image. Subsequently, the d -d space is indexed using a multidimensional index structure and queries are answered using the reduced dimensional index (see [17, 27] for details).

Clustering Clustering algorithms have been studied recently in the data mining domain (e.g., BIRCH, CLARANS, DBSCAN and CURE algorithms) [45, 35, 24, 14, 29]. The algorithms most related to this paper are those that discover patterns in low dimensional subspaces [1, 2]. In [1], Agarwal et. al. present an algorithm, called CLIQUE, to discover “dense” regions in all subspaces of the original data space. The algorithm works from lower to higher dimensionality subspaces: it starts by discovering 1-d dense units and iteratively discovers all dense units in each k -d subspace by building from the dense units in $(k-1)$ -d subspaces. In [2], Aggarwal et. al. present an algorithm, called PROCLUS, that clusters the data based on their correlation i.e. partitions the data into disjoint groups of correlated points. The authors use the hill climbing technique, popular in spatial cluster analysis, to determine the projected clusters. Neither CLIQUE, nor PROCLUS can be used as an LDR technique since they cannot discover clusters when the principal components are arbitrarily oriented. They can discover only those clusters that are correlated along one or more of the original dimensions. The above techniques are meant for discovering interesting patterns in the data; since correlation along arbitrarily oriented components is usually not that interesting to the user, they do not attempt to discover such correlation. On the contrary, the goal of LDR is efficient indexing; it must be able to discover such correlation in order to minimize the loss of information and make indexing efficient. Also, since the motivation of their work is pattern discovery and not indexing, they do not address the indexing and query processing issues which we have addressed in this paper. To the best of our knowledge, this is the first paper that proposes to exploit the local correlations in data for the purpose of indexing.

3 Identifying Correlated Clusters

In this section, we formally define the notion of correlated clusters and present an algorithm to discover such clusters in the data.

3.1 Definitions

In developing the algorithm to identify the correlated clusters, we will need the following definitions.

Definition 1 (Cluster and Subspace) Given a set \mathcal{A} of N points in a D -dimensional feature space, we define a *cluster* S as a set \mathcal{A}_S ($\mathcal{A}_S \subseteq \mathcal{A}$) of locally correlated points. Each cluster S is defined by $S = \langle \Phi_S, d_S, C_S, \mathcal{A}_S \rangle$ where:

Symbols	Definitions
N	Number of objects in the database
M	Maximum number of clusters desired
K	Actual number of clusters found ($K \leq M$)
D	Dimensionality of the original feature space
S_i	The i th cluster
C_i	Centroid of S_i
n_i	Size of S_i (number of objects)
\mathcal{A}_i	Set of points in S_i
Φ_i	The principal components of S_i
$\Phi_i^{(j)}$	The j th principal component of S_i
d_i	Subspace dimensionality of S_i
ϵ	Neighborhood range
$MaxReconDist$	Maximum Reconstruction distance
$FracOutliers$	Permissible fraction of outliers
$MinSize$	Minimum Size of a cluster
$MaxDim$	Maximum subspace dimensionality of a cluster
\mathcal{O}	Set of outliers

Table 1: Summary of symbols and definitions

- Φ_S are the principal components of the cluster, $\Phi_S^{(i)}$ denoting the i th principal component.
- d_S is the reduced dimensionality i.e. the number of dimensions retained. Obviously, the retained dimensions correspond to the first d_S principal components $\Phi_S^{(i)}, 1 \leq i \leq d_S$ while the eliminated dimensions correspond to the next $(D - d_S)$ components. Hence we use the terms (principal) components and dimensions interchangeably in the context of the transformed space.
- $C_S = [C_S^{(d_S+1)} \dots C_S^{(D)}]$ is the centroid, that stores, for each eliminated dimension $\Phi_i, (d_S + 1) \leq i \leq D$, a single constant which is “representative” of the position of every point in the cluster along this unrepresented dimension (as we are not storing their unique positions along these dimensions).
- \mathcal{A}_S is the set of points in the cluster

The reduced dimensionality space defined by $\Phi_S^{(i)}, 1 \leq i \leq d_S$ is called the *subspace* of S . d_S is called the subspace dimensionality of S . ■

Definition 2 (Reconstruction Vector) Given a cluster $S = \langle \Phi_S, d_S, C_S, \mathcal{A}_S \rangle$, we define the *reconstruction vector* $\overline{ReconVect}(Q, S)$ of a point Q from S as follows:

$$\overline{ReconVect}(Q, S) = \overline{\sum_{i=(d_S+1)}^D} (Q \bullet \Phi_S^{(i)} - C_S^{(i)}) \Phi_S^{(i)} \quad (1)$$

where $\overline{\sum}$ denotes vector addition and \bullet denotes scalar product (i.e. $Q \bullet \Phi_S^{(i)}$ is the projection of Q on $\Phi_S^{(i)}$ as shown in Figure 2). $(Q \bullet \Phi_S^{(i)} - C_S^{(i)})$ is the (scalar) distance of Q from the centroid along each eliminated dimension and $\overline{ReconVect}(Q, S)$ is the vector of these distances. ■

Definition 3 (Reconstruction Distance) Given a cluster $S = \langle \Phi_S, d_S, C_S, \mathcal{A}_S \rangle$, we now define the *reconstruction distance* (scalar) $ReconDist(Q, S, \mathcal{D})$ of a point Q from S . \mathcal{D} is the distance function used to

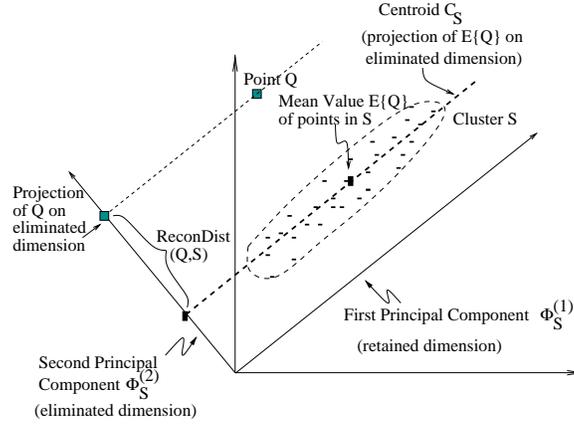


Figure 2: Centroid and Reconstruction Distance.

define the similarity between points in the HDFS. Let \mathcal{D} be an L_p metric i.e. $\mathcal{D}(P, P') = \|P - P'\|_p = [\sum_{i=1}^d (|P[i] - P'[i]|)^p]^{1/p}$. We define $ReconDist(Q, S, \mathcal{D})$ ³ as follows:

$$ReconDist(Q, S, \mathcal{D}) = ReconDist(Q, S, L_p) = \| \overline{ReconVect}(Q, S) \|_p = [\sum_{i=(d_S+1)}^D (|Q \bullet \Phi_S^{(i)} - C_S^{(i)}|)^p]^{1/p} \quad (2)$$

For any point Q mapped to the d_S -dimensional subspace of S , $\overline{ReconVect}(Q, S)$ represents the error in the representation i.e. the vector difference between the exact D -dimensional representation of Q and its approximate representation in the d_S -dimensional subspace of S . Higher the error, more the amount of distance information lost. When averaged over all points in S , we get the average information loss in S which is called the reconstruction error of S .

Definition 4 (Reconstruction Error) The reconstruction error $\bar{\epsilon}^2(S)$ of cluster S is defined as the mean square magnitude of $\overline{ReconVect}(Q, S)$ where $Q \in \mathcal{A}_S$:

$$\bar{\epsilon}^2(S) = E\{(\|\overline{ReconVect}(Q, S)\|_2)^2\} = \sum_{i=(d_S+1)}^D E\{(C_S^{(i)} - Q \bullet \Phi_S^{(i)})^2\} \quad (3)$$

where $E(X)$ denotes expected value of X .

3.2 Constraints on Correlated Clusters

Our objective in defining clusters is to identify low dimensional subspaces, one for each cluster, that can be indexed separately. We desire each subspace to have as low dimensionality as possible without losing too much distance information. In order to achieve the desired goal, each cluster must satisfy the following constraints:

1. **Reconstruction Distance Bound:** In order to restrict the maximum representation error of any point in the low dimensional subspace, we enforce the reconstruction distance of any point $P \in \mathcal{A}_S$ to satisfy the following condition: $ReconDist(P, S) \leq MaxReconDist$ where $MaxReconDist$ is a parameter specified by the user. This condition restricts the amount of information lost within each cluster and hence guarantees a high precision which in turn implies lower query cost.

³Assuming that \mathcal{D} is a fixed L_p metric, we usually omit the \mathcal{D} in $ReconDist(Q, S, \mathcal{D})$ for simplicity of notation.

2. **Dimensionality Bound:** For efficient indexing, we want the subspace dimensionality to be as low as possible while still maintaining high query precision. A cluster must not retain any more dimensions than necessary. In other words, it must retain the minimum number of dimensions required to accommodate the points in the dataset. Note that a cluster S can accommodate a point P only if $ReconDist(P, S) \leq MaxReconDist$. To ensure that the subspace dimensionality d_S is below the critical dimensionality of the multidimensional index structure (i.e. the dimensionality above which a sequential scan is better), we enforce the following condition: $d_S \leq MaxDim$ where $MaxDim$ is specified by the user.
3. **Choice of Centroid:** For each cluster S , we use PCA to determine the subspace i.e. Φ_S is the set of eigenvectors of the covariance matrix of \mathcal{A}_S sorted based on their eigenvalues. [20] shows that for a given choice of reduced dimensionality d_S , the reconstruction error $\bar{\epsilon}^2(S)$ is minimized by choosing the first d_S components among Φ_S and choosing C_S to be the mean value of the points (i.e. the centroid) projected on the eliminated dimensions. To minimize the information loss, we choose $C_S^{(i)} = E\{P \bullet \Phi_S^{(i)}\} = E\{P\} \bullet \Phi_S^{(i)}$ (see Figure 2).
4. **Size Bound:** Finally, we desire each cluster to have a minimum cardinality (number of points) : $n_S \geq MinSize$ where $MinSize$ is user-specified. The clusters that are too small are considered to be outliers.

The goal of the LDR algorithm described below is to discover the set $\mathcal{S} = S_1, S_2, \dots, S_K$ of K clusters (where $K \leq M$, M being the maximum number of clusters desired) that exists in the data and that satisfy the above constraints. The remaining points, that do not belong to any of the clusters, are placed in the outlier set \mathcal{O} .

3.3 The Clustering Algorithm

Since the LDR algorithm needs to perform *local* correlation analysis (i.e. PCA on subsets of points in the dataset rather than the whole dataset), we need to first identify the right subsets to perform the analysis on. This poses a cyclic problem: how do we identify the right subsets without doing the correlation analysis and how do we do the analysis without knowing the subsets. We break the cycle by using *spatial clusters* as an initial guess of the right subsets. Then we perform PCA on each spatial cluster individually. Finally, we ‘recluster’ the points based on the correlation information (i.e. principal components) to obtain the correlated clusters. The clustering algorithm is shown in Table 2. It takes a set of points \mathcal{A} and a set of clusters \mathcal{S} as input. When it is invoked for the first time, \mathcal{A} is the entire dataset and each cluster in \mathcal{S} is marked ‘empty’. At the end, each identified cluster is marked ‘complete’ indicating a completely constructed cluster (no further change); the remaining clusters remain marked ‘empty’. The points that do not belong to any of the clusters are placed to the outlier set \mathcal{O} . The details of each step is described below:

- **Construct Spatial Clusters**(Steps FC1 and FC2): The algorithm starts by constructing M spatial clusters where M is the maximum number of clusters desired. We use a simple single-pass partitioning-based spatial clustering algorithm to determine the spatial clusters [29, 35]. We first choose a set of $\mathcal{C} \subset \mathcal{A}$ of *well-scattered* points as the centroids such that points that belong to the same spatial cluster are not chosen to serve as centroids to different clusters. Such a set \mathcal{C} is called a *piercing* set [2]. We achieve this by ensuring that each point $P \in \mathcal{C}$ in the set is sufficiently far from any already chosen point $P' \in \mathcal{C}$ i.e. $Dist(P, P') > threshold$ for a user-defined threshold.⁴ This technique, proposed by Gonzalez [22], is guaranteed to return a piercing if no outliers are present. To avoid scanning though the whole database

⁴For subsequent invocations of FindClusters procedure during the iterative algorithm (Step 2 in Table 3), there may exist already completed clusters (does not exist during the initial invocation). Hence P must also be sufficiently far from all complete clusters formed so far i.e. $ReconDist(P, S) > threshold$ for each complete cluster S .

Clustering Algorithm
Input: Set of Points \mathcal{A} , Set of clusters \mathcal{S} (each cluster is either empty or complete)
Output: Some empty clusters are completed, the remaining points form the set of outliers \mathcal{O}
FindClusters ($\mathcal{A}, \mathcal{S}, \mathcal{O}$)
FC1: For each empty cluster, select a random point $P \in \mathcal{A}$ such that P is sufficiently far from all completed and valid clusters. If found, make P the centroid C_i and mark S_i valid.
FC2: For each point $P \in \mathcal{A}$, add P to the closest valid cluster S_i (i.e. $i = \operatorname{argmin}(\operatorname{Distance}(P, C_i))$) if P lies in the ϵ -neighborhood of C_i i.e. $\operatorname{Distance}(P, C_i) \leq \epsilon$.
FC3: For each valid cluster S_i , compute the principal components Φ_i using PCA. Remove all points from \mathcal{A}_i .
FC4: For each point $P \in \mathcal{A}$, find the valid cluster S_i that, among all the valid clusters requires the minimum subspace dimensionality $LD(P)$ to satisfy $\operatorname{ReconDist}(P, S_i) \leq \operatorname{MaxReconDist}$ (break ties arbitrarily). If $LD(P) \leq \operatorname{MaxDim}$, increment $V_i[j]$ for $j = 0$ to $(LD(P) - 1)$ and n_i .
FC5: For each valid cluster S_i , compute the subspace dimensionality d_i as: $d_i = \{j F_i[j] \leq \operatorname{FracOutliers} \text{ and } F_i[j - 1] > \operatorname{FracOutliers}\}$ where $F_i[j] = \frac{V_i[j]}{n_i}$.
FC6: For each point $P \in \mathcal{A}$, add P to the first valid cluster S_i such that $\operatorname{ReconDist}(P, S_i) \leq \operatorname{MaxReconDist}$. If no such S_i exists, add P to \mathcal{O} .
FC7: If a valid cluster S_i violates the size constraint i.e. $(\mathcal{A}_i < \operatorname{MinSize})$, mark it empty. Remove each point $P \in \mathcal{A}_i$ from S_i and add it to the first succeeding cluster S_j that satisfies $\operatorname{ReconDist}(P, S_j) \leq \operatorname{MaxReconDist}$ or to \mathcal{O} if there is no such cluster. Mark the other valid clusters complete. For each complete cluster S_i , map each point $P \in \mathcal{A}_i$ to the subspace and store it along with $\operatorname{ReconDist}(P, S, \mathcal{D})$.

Table 2: Clustering Algorithm

to choose the centroids, we first construct a random sample of the dataset and choose the centroids from the sample [2, 19, 24]. We choose the sample to be large enough (using Chernoff bounds [34]) such that the probability of missing clusters due to sampling is low i.e. there is at least one point from each cluster present in the sample with a high probability [24]. Once the centroids are chosen, we group each point $P \in \mathcal{A}$ with the closest centroid $C_{closest}$ if $\operatorname{Distance}(P, C_{closest}) \leq \epsilon$ and update the centroid to reflect the mean position of its group. If $\operatorname{Distance}(P, C_{closest}) > \epsilon$, we ignore P . The restriction of the neighborhood range to ϵ makes the correlation analysis *localized*. Smaller the value of ϵ , the more localized the analysis. At the same time, ϵ has to be large enough so that we get a sufficiently large number of points in the cluster which is necessary for the correlation analysis to be robust.

- **Compute PCs**(Step FC3): Once we have the spatial clusters, we perform PCA on each spatial cluster S_i individually to obtain the principal components $\Phi_S^{(i)}$, $i = [1, D]$ (see Appendix A for details on PCA). We do not eliminate any components yet. We compute the mean value M_i of the points in S_i so that we can compute $\operatorname{ReconDist}(P, S_i)$ in Steps FC4 and FC5 for any choice of subspace dimensionality d_i . Finally, we remove the points from the spatial clusters so that they can be reclustered as described in Step FC6.
- **Determine Subspace Dimensionality**(Steps FC4 and FC5): For each cluster S_i , we must retain no more dimensions than necessary to accommodate the points in the dataset (except the outliers). To determine the

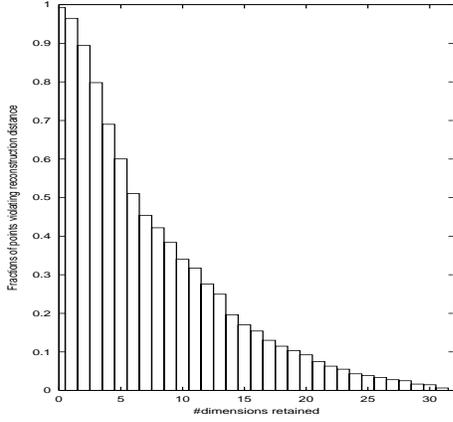


Figure 3: Determining subspace dimensionality (MaxDim=32).

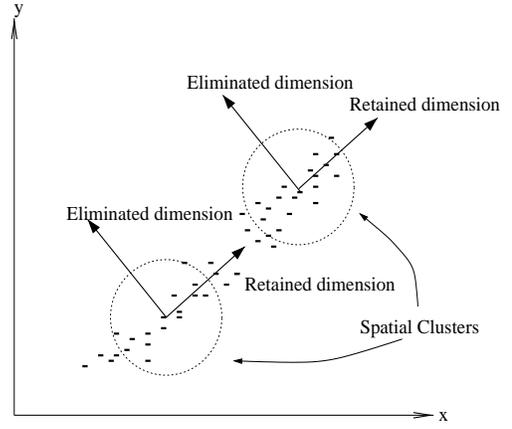


Figure 4: Splitting of correlated clusters due to initial spatial clustering.

number of dimensions d_i to be retained for each cluster S_i , we first determine, for each point $P \in \mathcal{A}$, the best cluster, if one exists, for placing P . Let $LD(P, S_i)$ denote the the least dimensionality needed for the cluster S_i to represent P with $ReconDist(P, S_i) \leq MaxReconDist$. Formally,

$$LD(P, S_i) = \{d \quad | \quad ReconDist(P, S_i) \leq MaxReconDist \text{ if } d_i \geq d \quad (4)$$

$$\text{and } ReconDist(P, S_i) > MaxReconDist \text{ otherwise} \quad (5)$$

In other words, the first $LD(P, S_i)$ PCs are just enough to satisfy the above constraint. Note that such a $LD(P, S_i)$ always exists for a non-negative $MaxReconDist$. Let $LD(P) = \min \{ LD(P, S_i) | S_i \text{ is a valid cluster} \}$. If $LD(P) \leq MaxDim$, there exists a cluster that can accommodate P without violating the dimensionality bound. Let $LD(P, S_i) = LD(P)$ (if there are multiple such clusters S_i , break ties arbitrarily). We say S_i is the “best” cluster for placing P since S_i is the cluster that, among all the valid clusters, needs to retain the minimum number of dimensions to accommodate P . P would satisfy the $ReconDist(P, S_i) \leq MaxReconDist$ bound if the subspace dimensionality d_i of S_i is such that $LD(P, S_i) \leq d_i \leq MaxDim$ and would violate it if $0 \leq d_i < LD(P, S_i)$. For each cluster S_i , we maintain this information as a count array $V_i[j], j = [0, MaxDim]$ where $V_i[j]$ is the number of points that, among the points chosen to be placed in S_i , would violate the $ReconDist(P, S_i) \leq MaxReconDist$ constraint if the subspace dimensionality d_i is j : so in this case (for point P), we must increment $V_i[j]$ for $j = 0$ to $(LD(P, S_i) - 1)$ and the total count n_i of points chosen to be placed in S_i . ($V_i[j]$ and n_i is initialized to 0 before FC4 begins). On the other hand, if $LD(P) > MaxDim$, there exists no cluster in which P can be placed without violating the dimensionality bound; so we do nothing.

At the end of the pass over the dataset, for each cluster S_i , we have computed $V_i[j], j = [0, MaxDim]$ and n_i . We use this to compute $F_i[j], j = [0, MaxDim]$ where $F_i[j]$ is the fraction of points that, among those chosen to be placed in S_i (during FC4), would violate the $ReconDist(P, S_i) \leq MaxReconDist$ constraint if the subspace dimensionality d_i is j i.e. $F_i[j] = \frac{V_i[j]}{n_i}$. An example of F_i from one of the experiments conducted on the real life dataset (cf. Section 5.3) is shown in Figure 3. We choose d_i to be as low as possible without too many points violating the reconstruction distance bound i.e. not more than $FracOutliers$ fraction of points in S_i where $FracOutliers$ is specified by the user. In other words, d_i is the minimum number of dimensions that must be retained so that the fraction of points that violate the $ReconDist(P, S_i) \leq MaxReconDist$ constraint is no more that $FracOutliers$ i.e. $d_i = \{j | F_i[j] \leq FracOutliers \text{ and } F_i[j - 1] > FracOutliers\}$. In Figure 3, d_i is 21 for $FracOutliers = 0.1$, 16 for

$FracOutliers = 0.2$ and 14 for $FracOutliers = 0.3$. We now have all the subspaces formed. In the next step, we assign the points to the clusters.

- **Recluster Points**(Step FC6): In the reclustering step, we reassign each point $P \in \mathcal{A}$ to a cluster S that covers P i.e. $ReconDist(P, S) \leq MaxReconDist$. If there exists no such cluster, P is added to the outlier set \mathcal{O} . If there exists just one cluster that covers P , P is assigned to that cluster. Now we consider the interesting case of multiple clusters covering P . In this case, there is a possibility that some of these clusters are actually parts of the same correlated cluster but has been split due to the initial spatial clustering. This is illustrated in Figure 4. Since points in a correlated cluster can be spatially distant from each other (e.g., form an elongated cluster in Figure 4) and spatial clustering only clusters spatially close points, it may end up putting correlated points in different spatial clusters, thus breaking up a single correlated cluster into two or more clusters. Although such ‘splitting’ does not affect the indexing cost of our technique for range queries and k-NN queries, it increases the cost of point search and deletion as multiple clusters may need to be searched in contrast to just one when there is no ‘splitting’. (cf. Section 4.2.1). Hence, we must detect these ‘broken’ clusters and merge them back together. We achieve this by maintaining the clusters in some fixed order (e.g., order in which they were created). For each point $P \in \mathcal{P}$, we check each cluster sequentially in that order and assign it to the first cluster that covers P . If two (or more) clusters are part of the same correlated cluster, most points will be covered by all of them but will *always* be assigned to only one of them, whichever appears first in the order. This effectively merges the clusters into one since only the first one will remain while the others will end up being almost empty and will be discarded due to the violation of size bound in FC7. Note that the $FracOutliers$ bound in Step FC5 still holds i.e. besides the points for which $LD(P) > MaxDim$, no more than $FracOutliers$ fraction of points can become outliers.
- **Map Points**(Step FC7): In the final step of the algorithm, we eliminate clusters that violate the size constraint. We remove each point from these clusters and add it to the first succeeding valid cluster S_j that satisfies the $ReconDist(P, S_j) \leq MaxReconDist$ bound or to \mathcal{O} otherwise. For the remaining clusters S_i , we map each point $P \in \mathcal{A}_i$ to the subspace by projecting P to $\Phi_i^{(j)}, 1 \leq j \leq d_i$ and refer it as the (d_i -d) image $Image(P, S_i)$ of P :

$$Image(P, S_i)[j] = P \bullet \Phi_i^{(j)} \text{ for } 1 \leq j \leq d_i \quad (6)$$

We refer to P as the (D -d) original $Original(Image(P, S_i), S_i)$ of its image $Image(P, S_i)$. We store the image of each point along with the reconstruction distance $ReconDist(P, S_i)$.

Since FindClusters chooses the initial centroids from a random sample, there is a risk of missing out some clusters. One way to reduce this risk is to choose a large number of initial centroids but at the cost of slowing down the clustering algorithm. We reduce the risk of missing clusters by trying to discover more clusters, if there exists, among the points returned as outliers by the initial invocation of FindClusters. We iterate the above process as long as new clusters are still being discovered as shown below:

Iterative Clustering	
(1)	FindClusters($\mathcal{A}, \mathcal{S}, \mathcal{O}$); /* initial invocation */
(2)	Let \mathcal{O}' be an empty set. Invoke FindClusters($\mathcal{O}, \mathcal{S}, \mathcal{O}'$). Make \mathcal{O}' the new outlier set i.e. $\mathcal{O} \leftarrow \mathcal{O}'$. If new clusters found, go to (2). Else return.

Table 3: Iterative Clustering Algorithm

The above iterative clustering algorithm is somewhat similar to the hill climbing technique, commonly used in spatial clustering algorithms (especially in partitioning-based clustering algorithms like k-means, k-medoids and CLARANS [29, 35]). In this technique, the “bad quality” clusters (the ones that violate the size bound) are discarded (Step FC7) and is replaced, if possible, by better quality clusters. However, unlike the hill climbing approach where all the points are reassigned to the clusters, we do not reassign the points already assigned to the ‘complete’ clusters. Alternatively, we can follow the hill climbing approach but it is computationally more expensive and requires more scans of the database [35].

Cost Analysis: We conclude this section with a analysis of the cost of the clustering algorithm. Let us first analyze the cost of the first invocation of the FindClusters procedure (where \mathcal{A} is the whole dataset). The centroid selection step (FC1) has a small cost since we are using a random sample and $|sample| \ll |\mathcal{A}|$. Step FC2 requires one pass through the dataset \mathcal{A} and has a time complexity of $O(NKD)$. Step FC3 has a complexity of $O(n_i D^2)$ for each cluster S_i and hence an overall complexity of $O(ND^2)$ (since $\sum_i n_i \leq N$). This step also has a memory requirement of $O(n_i D)$ for each cluster and hence a maximum of $O(max_i(n_i)D)$ which is smaller than the memory requirement of $O(ND)$ of GDR. This is an advantage of LDR over GDR: while the latter requires the whole dataset to fit in memory, the former requires only the points in the cluster to fit in memory. In either case, if the memory is too small, we can perform SVD on a sample rather than the whole data [27]. Step FC4 requires another pass through the database and has a time complexity of $O(ND^2K)$ (assuming $MaxDim$ is a constant). Step FC5 is a simple step with a complexity of $O(KD)$. Step FC6 requires a final pass through the database and has a time complexity of $O(ND^2K)$. Also, the first invocation of FindClusters accounts for most of the cost of the algorithm since the later invocations have much smaller sets as input and hence much smaller cost. Thus, the algorithm requires three passes through the dataset (FC2, FC4 and FC6) and a time complexity of $O(ND^2K)$.

4 Indexing Correlated Clusters

Having developed the technique to find the correlated clusters, we now shift our attention to how to use them for indexing. Our objective is to develop a data structure that exploits the correlated clusters to efficiently support range and k-NN queries over HDFSSs. The developed data structure must also be able to handle insertions and deletions.

4.1 Data Structure

The data structure, referred to as the global index structure (GI) (i.e. index on entire dataset), consists of separate multidimensional indices for each cluster, connected to a single root node. The global index structure is shown in Figure 5. We explain the various components in details below:

- *The Root Node R* of GI contains the following information for each cluster S_i : (1) a pointer to the root node R_i (i.e. the address of disk block containing R_i) of the cluster index I_i (the multidimensional index on S_i), (2) the principal components Φ_i (3) the subspace dimensionality d_i and (4) the centroid C_i . It also contains an access pointer O to the outlier cluster \mathcal{O} . If there is an index on \mathcal{O} (discussed later), O points to the root node of that index; otherwise, it points to the start of the set of blocks on which the outlier set resides on disk. R may occupy one or more disk blocks depending on the number of clusters K and original dimensionality D .
- *The Cluster Indices:* We maintain a multidimensional index I_i for each cluster S_i in which we store the reduced dimensional representation of the points in S_i . However, instead of building the index I_i on the d_i -d subspace of S_i defined by $\Phi_i^{(j)}, 1 \leq j \leq d_i$, we build I_i on the $(d_i + 1)$ -d space, the first d_i dimensions of which are defined by $\Phi_i^{(j)}, 1 \leq j \leq d_i$ as above while the $(d_i + 1)$ th dimension is defined by the

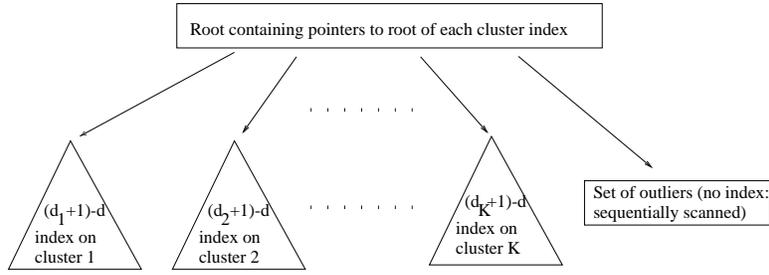


Figure 5: The global index structure

reconstruction distance $ReconDist(P, S_i, \mathcal{D})$. Including reconstruction distance as a dimension helps to improve query precision (as explained later). We redefine the image $NewImage(P, S_i)$ of a point $P \in \mathcal{A}_i$ as a $(d_i + 1)$ -d point (rather than a d_i -d point), incorporating the reconstruction distance as the $(d_i + 1)$ th dimension:

$$NewImage(P, S_i)[j] = Image(P, S_i)[j] = P \bullet \Phi_i^{(j)} \text{ for } 1 \leq j \leq d_i \quad (7)$$

$$= ReconDist(P, S_i, \mathcal{D}) \text{ for } j = d_i + 1 \quad (8)$$

The $(d_i + 1)$ -d cluster index I_i is constructed by inserting the $(d_i + 1)$ -d images (i.e. $NewImage(P, S_i)$) of each point $P \in \mathcal{A}_i$ into the multidimensional index structure using the insertion algorithm of the index structure. Any disk-based multidimensional index structure (e.g., R-tree [25], X-tree [5], M-tree [11], Hybrid Tree [7]) can be used for this purpose. We used the hybrid tree in our experiments since it is a space partitioning index structure (i.e. has “dimensionality-independent” fanout), is more scalable to high dimensionalities in terms of query cost and can support arbitrary distance metrics [7, 38, 9].

- *The Outlier Index:* For the outlier set \mathcal{O} , we may or may not build an index depending on whether the original dimensionality D is below or above the critical dimensionality. In this paper, we assume that D is above the critical dimensionality of the index structure and hence choose not to index the outlier set (i.e. use sequential scan for it).

Like other database index trees (e.g., B-tree, R-tree), the global index (GI) shown in Figure 5 is disk-based. But it may not be perfectly height balanced i.e. all paths from R to leaf may not be of exactly equal length. The reason is that the sizes and the dimensionalities may differ from one cluster to another causing the cluster indices to have different heights. We found that GI is *almost* height balanced (i.e. the difference in the lengths of *any* two paths from R to leaf is never more than 1 or 2) due to the size bound on the clusters (see Appendix D for details). Also, its height cannot exceed the height of the original space index by more than 1 (see Appendix D for details).

To guarantee the correctness of our query algorithms (i.e. to ensure no false dismissals), we need to show that the cluster index distances *lower bounds* the actual distances in the original D -d space [17, 16]. In other words, for any two D -d points P and Q , $\mathcal{D}(NewImage(P, S_i), NewImage(Q, S_i))$ must always lower bound $\mathcal{D}(P, Q)$.

Lemma 1 (Lower Bounding Lemma) $\mathcal{D}(NewImage(P, S_i), NewImage(Q, S_i))$ always lower bounds $\mathcal{D}(P, Q)$. (Proof in Appendix B).

Note that instead of incorporating reconstruction distance as the $(d_i + 1)$ th dimension, we could have simply constructed GI with each cluster index I_i defined on the corresponding d_i -d subspace $\Phi_i^{(j)}$, $1 \leq j \leq d_i$. Since the lower bounding lemma holds for the d_i -d subspaces (as shown in [17]), the query processing algorithms described below would have been correct. The reason we use $(d_i + 1)$ -d subspace is that the distances in the

$(d_i + 1)$ -d subspace upper bounds the distances in the d_i -d subspace and hence provides a tighter lower bound to distances in the original D -d space:

$$\begin{aligned} \mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i)) &= \\ & [\mathcal{D}(\text{Image}(P, S_i), \text{Image}(Q, S_i))^p + |(\text{ReconDist}(P, S_i, \mathcal{D}) - \text{ReconDist}(Q, S_i, \mathcal{D}))|^p]^{1/p} \\ & \Rightarrow \mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i)) \geq \mathcal{D}(\text{Image}(P, S_i), \text{Image}(Q, S_i)) \quad (9) \end{aligned}$$

Furthermore, the difference between the two (i.e. $\mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i))$ and $\mathcal{D}(\text{Image}(P, S_i), \text{Image}(Q, S_i))$) is usually significant when computing the distance of the query from a point in the cluster: Say, P is a point in S_i and Q is the query point. Due to the reconstruction distance bound, $\text{ReconDist}(P, S_i, \mathcal{D})$ is *always* a small number ($\leq \text{MaxReconDist}$). On the other hand, $\text{ReconDist}(Q, S_i, \mathcal{D})$ can have any arbitrary value and is usually much larger than $\text{ReconDist}(P, S_i, \mathcal{D})$, thus making the difference quite significant. This makes the distance computations in the $(d_i + 1)$ -d more optimistic than that in the d_i -d index and hence a better estimate of the distances in the original D -d space. For example, for a range query, the range condition ($\mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i)) \leq \rho$) is more optimistic (i.e. satisfies fewer objects) than the range condition ($\mathcal{D}(\text{Image}(P, S_i), \text{Image}(Q, S_i)) \leq \rho$), leading to fewer false positives. The same is true for k-NN queries. Fewer false positives imply lower query cost. At the same time, adding a new dimension also increases the cost of the query. Our experiments show that decrease in the query cost from fewer false positives offsets the increase of the cost of the adding a dimension, reducing the overall cost of the query significantly (cf. Section 5, Figure 12).

4.2 Query Processing over the Global Index

In this section, we discuss how to execute similarity queries efficiently using the index structure described above (cf. Figure 5). We describe the query processing algorithm for point, range and k-NN queries. For correctness, the query processing algorithm must guarantee that it always returns exactly the same answer as the query on the original space [17, 16]. Often dimensionality reduction techniques do not satisfy the correctness criteria [27, 42]. We show that all our query processing algorithms satisfy the above criteria.

4.2.1 Point Search

To find an object O , we first find the cluster that contains O . It is the first cluster S (in the order mentioned in Step FC6) for which the reconstruction distance bound is satisfied. If such a cluster S exists, we compute $\text{NewImage}(O, S)$ and find it in the corresponding index by invoking the point search algorithm of the index structure. The point search returns the object if it exists in the cluster, otherwise it returns null. If no such cluster S exists, O must be, if at all, in \mathcal{O} . So we sequentially search through \mathcal{O} and return it if it exists in \mathcal{O} .

4.2.2 Range Queries

A range query $\mathcal{Q} = \langle Q, \rho, \mathcal{D} \rangle$ retrieves all objects O in the database that satisfies the range condition $\mathcal{D}(Q, O) \leq \rho$. The algorithm proceeds as follows (see Appendix C for pseudocode). For each cluster S_i , we map the query anchor Q to its $(d_i + 1)$ -d image Q_i (using the principal components Φ_i and subspace dimensionality d_i stored in the root node R of GI) and execute a range query (with the same range ρ) on the corresponding cluster index I_i by invoking the procedure `RangeSearchOnClusterIndex` on the root node R_i of I_i . `RangeSearchOnClusterIndex` is the standard R-tree-style recursive range search procedure that starts from the root node and explores the tree in a depth-first fashion. It examines the current node T : if T is a non-leaf node, it recursively searches each child

node N of T that satisfies the condition $MINDIST(Q, N, \mathcal{D}) \leq \rho$ (where $MINDIST(Q, N, \mathcal{D})$ denotes the minimum distance of the $(d_i + 1)$ -d image of query point to the $(d_i + 1)$ -d bounding rectangle of N based on distance function \mathcal{D} [26, 39]); if T is a leaf node, it retrieves each data item O stored in T (which is the *NewImage* of the original D -d object) that satisfies the range condition $\mathcal{D}(Q, O) \leq \rho$ in the $(d_i + 1)$ -d space, accesses the full D -dimensional tuple on disk to determine whether it is a false positive and adds it to the result set if it is not a false positive (i.e. it also satisfies the range condition $\mathcal{D}(Q, O) \leq \rho$ in the original D -d space). After all the cluster indices are searched, we add all the qualifying points from among the outliers to the result by performing a sequential scan on \mathcal{O} . Since the distance in the index space lower bounds the distance in the original space (cf. Lemma 1), the above algorithm cannot have any false dismissals. The algorithm cannot have any false positives either as they are filtered out before adding to the result set. The above algorithm thus returns exactly the same answer as the query on the original space.

In the above discussion, we assumed that we store the reduced representation of the points (i.e. the ‘NewImage’s) in the leaf pages of the cluster indices. Another option was to store the original D -d point in the leaf pages (although the index is built on the reduced space). With the former option, the index will have much fewer leaf nodes than the latter due to the smaller representation. On the other hand, in the latter case, the false positives can be eliminated at the leaf page level while the former would require an additional page access into the relation (where the full tuple is stored) to eliminate false positives. Since the index is usually a secondary index, we assume that for each match, we need to access the full tuple anyway (to retrieve the additional attributes). In that case, the extra cost of the former option is that of additional page accesses for *only* the false positives (see Section 5.1 for the details on the cost computations). Our experiments show that our technique usually operates in a high precision zone ($> 90\%$) i.e. has very few false positives. The experiments also show that the smaller size of the indices in the former approach saves enough query cost to compensate the few extra I/Os due to false positives. Hence we store just the *NewImages* in the leaf pages of the index structure.

4.2.3 k Nearest Neighbor Queries

A k-NN query $Q = \langle Q, k, \mathcal{D} \rangle$ retrieves a set \mathcal{R} of k objects such that for any two objects $O \in \mathcal{R}, O' \notin \mathcal{R}$, $\mathcal{D}(Q, O) \leq \mathcal{D}(Q, O')$. The algorithm for k-NN queries is shown in Table 4. Like the basic k-NN algorithm [26, 39], the algorithm uses a priority queue *queue* to navigate the nodes/objects in the database in increasing order of their distances from Q . Note that we use a single queue to navigate the entire global index i.e. we explore the nodes/objects of all the cluster indices in an intermixed fashion and do not require separate queues to navigate the different clusters. Each entry in *queue* is either a node or an object and stores 3 fields: the id of the node/object T it corresponds to, the cluster S it belongs to and its distance *dist* from the query anchor Q . The items (i.e. nodes/objects) are prioritized based on *dist* i.e. the smallest item appears at the top of the queue (min-priority queue). For nodes, the distance is defined by *MINDIST* while for objects, it is the point-to-point distance [26, 39]. Initially, for each cluster, we map the query anchor Q to its $(d_i + 1)$ -d image Q_i using the information stored in the root node R of GI (Line 2). Then, for each cluster index I_i , we compute the distance $MINDIST(Q_i, R_i, \mathcal{D})$ of Q_i from the root node R_i of I_i and push R_i into *queue* along with the distance and the id of the cluster S_i to which it belongs (Line 3). We also fill the set *temp* with the k closest neighbors of Q among the outliers by sequentially scanning through \mathcal{O} (Line 4).

After these initialization steps, we start navigating the index by popping the item from the top of *queue* at each step (Line 11). If the popped item is an object, we compute the distance of the original D -d object (by accessing the full tuple on disk) from Q and append it to *temp* (Lines 12-14). If it a node, we compute the distance of each of its children to the appropriate query image $Q_{top.S}$ (where *top.S* denotes the cluster which *top* belongs to) and push them into the queue (Lines 15-20). Note that the image for each cluster is computed just

k-NNSearch(Query $Q = Q, k, \mathcal{D}$)	
1	for (i=1; i \leq K; i++)
2	$Q_{S_i} \leftarrow \text{NewImage}(Q, S_i)$;
3	$\text{queue.push}(S_i, R_i, \text{MINDIST}(Q_i, R_i, \mathcal{D}))$;
4	Add to temp the k closest neighbors of Q among \mathcal{O} (using linear scan)
5	while (not $\text{queue.IsEmpty}()$)
6	$\text{top} = \text{queue.Top}()$;
7	for each object O in temp such that $O.\text{dist} \leq \text{top}.\text{dist}$
8	$\text{temp} \leftarrow \text{temp} - O$;
9	$\text{result} = \text{result} \cup O$;
10	$\text{retrieved}++$;
11	if ($\text{retrieved} = k$) return result ;
12	$\text{queue.Pop}()$;
13	if top.T is an object
14	$\text{top}.\text{dist} = \mathcal{D}(Q, \text{Original}(\text{top.T}, \text{top.S}))$;
15	$\text{temp} = \text{temp} \cup \text{top.T}$;
16	else if top.T is a leaf node
17	for each object O in top.T
18	$\text{queue.push}(\text{top.S}, O, \mathcal{D}(Q_{\text{top.S}}, O))$;
19	else /* top.T is an index node */
20	for each child N of top.T
21	$\text{queue.push}(\text{top.S}, N, \text{MINDIST}(Q_{\text{top.S}}, N, \mathcal{D}))$;

Table 4: k-NN Query.

once (in Step 2) and is reused here. We move an object O from temp to result only when we are sure that it is among the k nearest neighbors of Q i.e. there exists no object $O' \notin \text{result}$ such that $\mathcal{D}(O', Q) < \mathcal{D}(O, Q)$ and $|\text{result}| < k$. The second condition is ensured by the exit condition in Line 11. The condition $O.\text{dist} \leq \text{top}.\text{dist}$ in Line 7 ensures that there exists no *unexplored* object O' such that $\mathcal{D}(O', Q) < \mathcal{D}(O, Q)$. The proof is simple: $O.\text{dist} \leq \text{top}.\text{dist}$ implies $O.\text{dist} \leq \mathcal{D}(\text{NewImage}(O', S), \text{NewImage}(Q, S))$ for any unexplored object O' in a cluster S (by the property of min-priority queue) which in turn implies $\mathcal{D}(O, Q) \leq \mathcal{D}(O', Q)$ (since $\mathcal{D}(\text{NewImage}(O', S), \text{NewImage}(Q, S))$ lower bounds $\mathcal{D}(O', Q)$, see Lemma 1). By inserting the objects in temp (i.e. already explored items) into result in increasing order of their distances in the original D-d space (by keeping temp sorted), we also ensure there exists no *explored* object O' such that $\mathcal{D}(O', Q) < \mathcal{D}(O, Q)$. This shows that the algorithm returns the correct answer i.e. the exact set of objects as the query in the original D-d space. It is also easy to show that the algorithm is I/O optimal.

Lemma 2 (Optimality of k-NN algorithm) *The k-NN algorithm is optimal i.e. it does not explore any object outside the range of k th nearest neighbor. (Proof in Appendix C).*

4.3 Modifications

We assume that the data is static in order to build the index. However, we must support subsequent insertions/deletions of the objects to/from the index efficiently. To insert an object O , we find the first cluster S (in the order mentioned earlier) for which the reconstruction distance bound is satisfied i.e. $\text{ReconDist}(O, S, \mathcal{D}) \leq$

ReconError. If such a cluster exists, we compute $NewImage(O, S)$ and insert it into the corresponding index using the insertion algorithm of the index structure. Otherwise, we append O to \mathcal{O} .

The deletion algorithm is also simple. To delete an object O , we first find O by invoking the point search algorithm (cf. Section 4.2.1). If it is found in a cluster, we delete it using the deletion algorithm of the index structure; else if it is found in \mathcal{O} , we delete it from \mathcal{O} ; else, we return not found.

If the database is dynamic (i.e. frequent insertions and deletions), the principal components need to be updated from time to time. One option is to repeat the entire clustering algorithm and construct the index structure from scratch. This can be done more efficiently using techniques proposed by Ravi Kanth et. al. [27]. The idea is to use aggregate data, obtained from the cluster indices, to recompute the principal components for each cluster and then incorporate the new components back into the cluster indices. [27] shows that this technique improves the recomputation time significantly without degrading the quality of the index structure. We can use their approach to handle dynamic databases. On the other hand, if the database is more or less static (i.e. insertions and deletions are rare) as is often the case [17, 15], such recomputations are not necessary.

5 Experiments

In this section, we present the results of an extensive empirical study we have conducted to (1) evaluate the effectiveness of LDR as a high dimensional indexing technique and (2) compare it with other techniques, namely, GDR, original space indexing (OSI) and linear scan. We conducted our experiments on both synthetic and real-life datasets. The major findings of our study can be summarized as follows:

- **High Precision:** LDR provides up to an order of magnitude improvement in precision over the GDR technique at the same reduced dimensionality. This indicates that LDR can achieve the same reduction as GDR with significantly lower loss of distance information.
- **Low Query Cost:** LDR consistently outperforms other indexing techniques, namely GDR, original space indexing and sequential scan, in terms of query cost (combined I/O and CPU costs) for both synthetic and real-life datasets.

Thus, our experimental results validate the thesis of this paper that LDR is an effective indexing technique for high dimensional datasets. All experiments reported in this section were conducted on a Sun Ultra Enterprise 450 machine with 1 GB of physical memory and several GB of secondary storage, running Solaris 2.5.

5.1 Experimental Methodology

We conduct the following two sets of experiments to evaluate the LDR technique and compare it with other indexing techniques.

Precision Experiments Due to dimensionality reduction, both GDR and LDR, cause loss of distance information (e.g., in Figure 15 in Appendix A, the distance between D and E is lost due to elimination of the second principal component). More the number of dimensions eliminated, more the amount of information lost. We measure this loss by *precision* defined as $Precision = \frac{|R_{original}|}{|R_{reduced}|}$ where $R_{reduced}$ and $R_{original}$ are the sets of answers returned by the range query on the reduced dimensional space and the original HDFS respectively [27]. For k -NN queries, $R_{original}$ is the set of k actual answers while $R_{reduced}$ is the set of objects we need to explore before being sure that we seen all the k actual answers. Note that the set $(R_{reduced} - R_{original})$ represent the false positives; so $Precision = \frac{1}{1 + \frac{|false\ positives|}{|R_{original}|}}$. We repeat that since our algorithms guarantee that the user always gets back the correct set $R_{original}$ of answers (as if the query executed in the original HDFS), precision does *not* measure the quality of the answers returned to the user but just the information loss incurred by the DR technique

and hence the query cost. For a DR technique, if we fix the reduced dimensionality, the higher the precision, the lower the cost of the query, the more efficient the technique. We compare the GDR and LDR techniques based on precision at fixed reduced dimensionalities.

Cost Experiments We conducted experiments to measure the query cost (I/O and CPU costs) for each of the following four indexing techniques. We describe how we compute the I/O and CPU costs of the techniques below.

- *Linear Scan*: In this technique, we perform a simple linear scan on the original high dimensional dataset. The I/O cost in terms of sequential disk accesses is $\frac{N*(D*sizeof(float)+sizeof(id))}{PageSize}$. Since $sizeof(id) \ll (D * sizeof(float))$, we will ignore the $sizeof(id)$ henceforth. Assuming sequential I/O is 10 times faster than random I/O, the cost in terms of the random accesses is $\frac{N*sizeof(float)*D}{10*PageSize}$. The CPU cost is the cost of computing the distance of the query from each point in the database.
- *Original Space Indexing (OSI)*: In this technique, we build the index on the original HDFS itself using a multidimensional index structure. We use the hybrid tree as the index structure. The I/O cost (in terms of random disk accesses) of the query is the number of nodes of the index structure accessed. The CPU cost is the CPU time (excluding I/O wait) required to navigate the index and return the answers.
- *GDR*: In this technique, we perform PCA on the original dataset, retain the first few principal components (depending on the desired reduced dimensionality) and index the reduced dimensional space using the hybrid tree index structure. In this case, the I/O cost has 2 components: index page accesses (discussed in OSI) and accessing the full tuples in the relation for false positive elimination (post processing cost). The post processing cost can be one I/O per false positives in the worst case. However, as observed in [23], this assumption is overly pessimistic (and is confirmed by our experiments). We, therefore, assume the postprocessing I/O cost to be $\frac{num_false_positives}{2}$. The total I/O cost (in number of random disk accesses) is $index_page_access_cost + \frac{num_false_positives}{2}$. The CPU cost is the sum of the index CPU cost and the post processing CPU cost i.e. cost of computing the distance of the query from each of the false positives.
- *LDR*: In this technique, we index each cluster using the hybrid tree multidimensional index structure and used a linear scan for the outlier set. For LDR, the I/O cost of a query has 3 components: index page accesses for each cluster index, linear scan on the outlier set and accessing the full tuples in the relation (post processing cost). The total index page access cost is the total number of nodes accessed of all the cluster indices combined. The number of sequential disk accesses for the outlier scan is $\frac{|O|*D*sizeof(float)}{PageSize}$. The cost of outlier scan in terms of random accesses is $\frac{|O|*sizeof(float)*D}{10*PageSize}$. The postprocessing I/O cost is $\frac{num_false_positives}{2}$ (as discussed above). The total I/O cost (in number of random disk accesses) is $index_page_access_cost + \frac{|O|*sizeof(float)*D}{10*PageSize} + \frac{num_false_positives}{2}$. Similarly, the CPU cost is the sum of the index CPU cost, outlier scan CPU cost (i.e. cost of computing the distance of the query from each of the outliers) and the post processing cost (i.e. cost of computing the distance of the query from each of the false positives).

We chose the hybrid tree as the index structure for our experiments since it is a space partitioning index structure (“dimensionality-independent” fanout) and has been shown to scale to high dimensionalities [7, 38, 9].

⁵ We use a page size of 4KB for all our experiments.

⁵The performance gap between our technique and the other techniques was even greater with SR-tree [28] as the index structure due to higher dimensionality curse [7]. We do not report those results here but can be found in the full version of the paper [8].

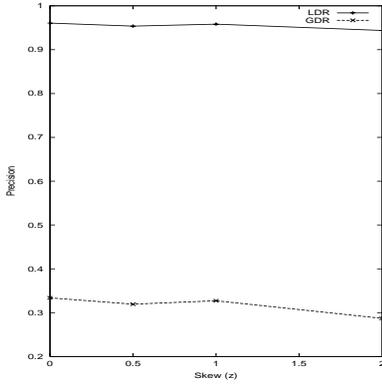


Figure 6: Sensitivity of precision to skew.

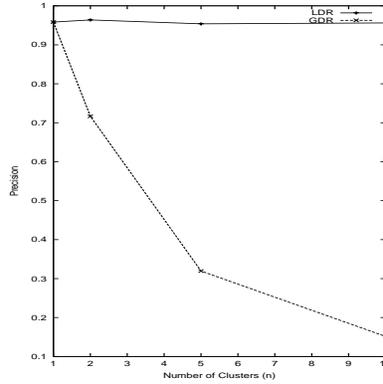


Figure 7: Sensitivity of precision to number of clusters.

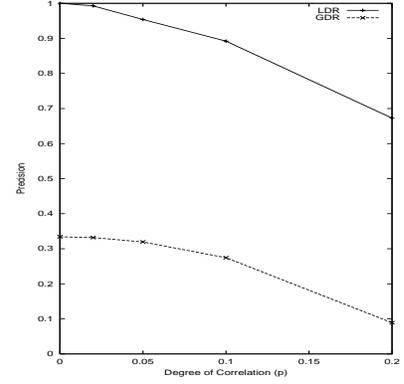


Figure 8: Sensitivity of precision to degree of correlation.

5.2 Experimental Results - Synthetic Data Sets

Synthetic Data Sets and Queries In order to generate the synthetic data, we use a method similar to that discussed in [45] but appropriately modified so that we can generate the different clusters in subspaces of different orientations and dimensionalities. The synthetic dataset generator is described in Appendix F. The input parameters to the data generator and their default values are shown in Table 6 (Appendix F).

We generated 100 range queries by selecting their query anchors randomly from the dataset and choosing a range value such that the average query selectivity is about 2%. We tested with only range queries since the k -NN algorithm, being optimal, is identical to the range query with the range equal to the distance of the k th nearest neighbor from the query (Lemma 3). We use L_2 distance (Euclidean) as the distance metric. All our measurements are averaged over the 100 queries.

Precision Experiments In our first set of experiments, we carry out a sensitivity analysis of the GDR and LDR techniques to parameters like skew in the size of the clusters (z_{size}), number of clusters (k) and degree of correlation (p). In each experiment, we vary the parameter of interest while the remaining parameters are fixed at their default values. We fix the reduced dimensionality of the GDR technique to 15. We fix the average subspace dimensionality of the clusters (i.e. $\sum_{i=1}^K \frac{n_i d_i}{K}$) also to 15 by choosing $FracOutliers$ and $MaxReconDist$ appropriately ($FracOutliers = 0.1$ and $MaxReconDist = 0.5$). Figure 6 compares the precision of the LDR technique with that of GDR for various value of z_{size} . LDR achieves about 3 times higher precision compared to GDR i.e. the latter has more than three times the number of false positives as the former. The precision of neither technique changes significantly with the skew. Figure 7 compares the precision of the two techniques for various values of k . As expected, for one cluster, the two techniques are identical. As k increases, the precision of GDR deteriorates while that of LDR is independent of the number of clusters. For $k = 10$, LDR is almost an order of magnitude better compared to GDR in terms of precision. Figure 8 compares the two techniques for various values of p . As the degree of correlation decreases (i.e. the value of p increases), the precision of both techniques drop but LDR outperforms GDR for all values p . Figure 9 shows the variation of the precision with the reduced dimensionality. For the GDR technique, we vary the reduced dimensionality from 15 to 60. For the LDR technique, we vary the $FracOutliers$ from 0.2 to 0.01 (0.2, 0.15, 0.1, 0.05, 0.02, 0.01) causing the average subspace dimensionality to vary from 7 to 42 (7, 10, 12, 14, 23 and 42) ($MaxDim$ was 64). The precision of both techniques increase with the increase in reduced dimensionality. Once again, LDR consistently outperforms GDR at all dimensionalities. The above experiments show that LDR is a more effective dimensionality reduction technique as it can achieve the same reduction as GDR with significantly lower loss of information (i.e. high

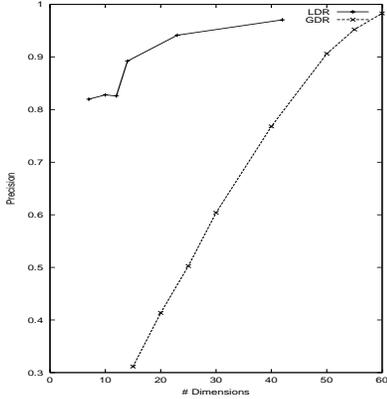


Figure 9: Sensitivity of precision to reduced dimensionality.

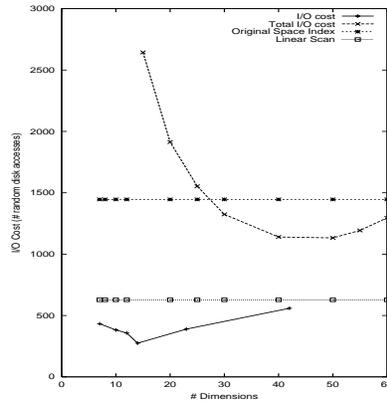


Figure 10: Comparison of LDR, GDR, Original Space Indexing and Linear Scan in terms of I/O cost. For linear scan, the cost is computed as: $\frac{\text{num_sequential_disk_accesses}}{10}$.

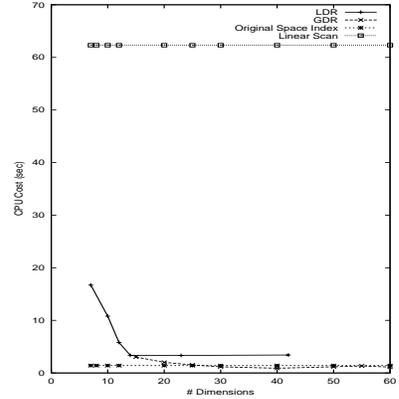


Figure 11: Comparison of LDR, GDR, Original Space Indexing and Linear Scan in terms of CPU cost.

precision) and hence significantly lower cost as confirmed in the cost experiments described next.

Cost Experiments We compare the 4 techniques, namely LDR, GDR, OSI and Linear Scan, in terms of query cost for the synthetic dataset. Figure 10 compares the I/O cost of the 4 techniques. Both the LDR and GDR techniques have U-shaped cost curves: when the reduced dimensionality is too low, there is a high degree of information loss leading to a large number of false positives and hence a high post-processing cost; when it is too high, the index page access cost becomes too high due to dimensionality curse. The optimum points lies somewhere in the middle: it is at dimensionality 14 (about 250 random disk accesses) for LDR and at 40 (about 1200 random disk accesses) for GDR. The I/O cost of OSI and Linear Scan is obviously independent of the reduced dimensionality. LDR significantly outperforms all the other 3 techniques in terms of I/O cost. The only technique that comes close to LDR in terms of I/O cost is the linear scan (but LDR is 2.5 times better as the latter performs 6274 sequential accesses \sim 627 random accesses). However, linear scan loses out mainly due to its high CPU cost shown in Figure 11. While LDR, GDR and OSI techniques have similar CPU cost (at their respective optimum points), the CPU cost linear scan is almost two orders of magnitude higher than the rest. LDR has slightly higher CPU cost compared to GDR and OSI since it uses linear scan for the outlier set: however, the savings in the I/O cost over GDR and OSI (by a factor of 5-6) far offsets the slightly higher CPU cost.

5.3 Experimental Results - Real-Life Data Sets

Description of Dataset Our real-life data set (COLHIST dataset [7]) comprises of 8×8 color histograms (64-d data) extracted from about 70,000 color images obtained from the Corel Database (<http://corel.digitalriver.com/>) and is available online at the UCI KDD Archive web site (<http://kdd.ics.uci.edu/databases/CorelFeatures>). We generated 100 range queries by selecting their query anchors randomly from the dataset and choosing a range value such that the average query selectivity is about 0.5%. All our measurements are averaged over the 100 queries.

Cost Experiments First, we evaluate the impact of adding *ReconDist* as an additional dimension of each cluster in the LDR technique. Figure 12 shows that the additional dimension reduces the cost of the query significantly. We performed the above experiment on the synthetic dataset as well and observed a similar result.

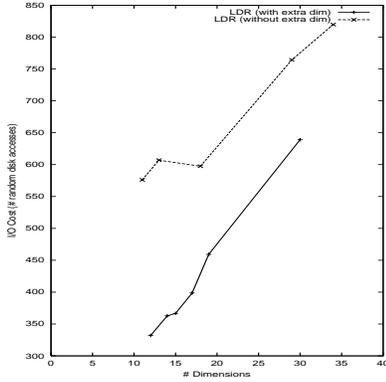


Figure 12: Effect of adding the extra dimension.

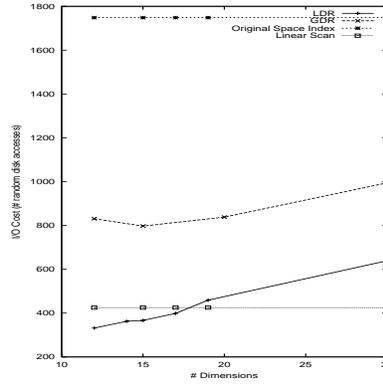


Figure 13: Comparison of LDR, GDR, Original Space Indexing and Linear Scan in terms of I/O cost. For linear scan, the cost is computed as: $\frac{\text{num_sequential_disk_accesses}}{10}$.

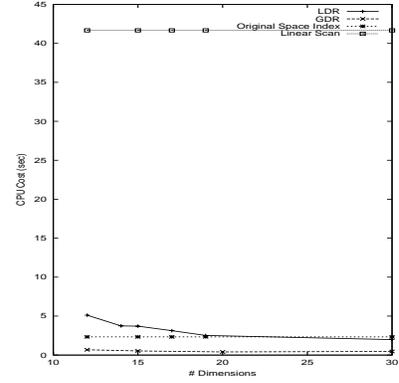


Figure 14: Comparison of LDR, GDR, Original Space Indexing and Linear Scan in terms of CPU cost.

⁶ Figure 13 compares the 4 techniques, namely LDR, GDR, OSI and Linear Scan, in terms of I/O cost. LDR outperforms all other techniques significantly. Again, the only technique that come close to LDR in I/O cost (i.e. number of random disk accesses) is the linear scan. However, again, linear scan turns out to significantly worse compared to LDR in terms of the overall cost due to its high CPU cost as shown in Figure 14.

6 Conclusion

With numerous emerging applications requiring efficient access to high dimensional datasets, there is a need for scalable techniques to indexing high dimensional data. In this paper, we proposed local dimensionality reduction (LDR) as an approach to indexing high dimensional spaces. We developed an algorithm to discover the locally correlated clusters in the dataset and perform dimensionality reduction on each of them individually. We presented an index structure that exploits the correlated clusters to efficiently support similarity queries over high dimensional datasets. We have shown that our query processing algorithms are correct and optimal. We conducted an extensive experimental study with synthetic as well as real-life datasets to evaluate the effectiveness of our technique and compare it to GDR, original space indexing and linear scan techniques. Our results demonstrate that our technique (1) reduces the dimensionality of the data with significantly lower loss in distance information compared to GDR, outperforming GDR by almost an order of magnitude in terms of query precision (for the same reduced dimensionality) and (2) significantly outperforms all the other 3 techniques (namely, GDR, original space indexing and linear scan) in terms of the query cost for both synthetic and real-life datasets.

7 Acknowledgements

We thank David Eppstein and Padhraic Smyth for the useful discussions on the clustering algorithm. We thank Kriengkrai Porkaew for the discussions and his help with the implementation. We thank Corel Corporation for making the large collection of images used in the COLHIST dataset available to us. Our PCA implementation is built on top of the Meschach Library downloaded from <http://www.netlib.org/c/meschach/>.

⁶We also analyzed the sensitivity of the LDR technique to the *MaxReconDist* parameter. The result is included in Appendix G.

References

- [1] R. Agarwal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic subspace clustering of high dimensional data for data mining applications. *Proc. of SIGMOD*, 1998.
- [2] C. Aggarwal, C. Procopiuc, J. Wolf, P. Yu, and J. Park. Fast algorithms for projected clustering. *Proc. of SIGMOD*, 1999.
- [3] M. Ankerst, M. Breunig, H. Kriegel, and J. Sander. Optics: ordering points to identify the clustering structure. *Proc. of SIGMOD*, 1999.
- [4] D. Barbara, W. DuMouchel, C. Faloutsos, P. Haas, J. Hellerstein, Y. Ionnidis, H. Jagadish, T. Johnson, R. Ng, V. Poosala, K. Ross, and K. Sevcik. The new jersey data reduction report. *Data Engineering*, 20(4), 1997.
- [5] S. Berchtold, D. A. Keim, and H. P. Kriegel. The x-tree: An index structure for high-dimensional data. *Proc. of VLDB*, 1996.
- [6] K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft. When is “nearest neighbor” meaningful? *Proc. of ICDDT*, 1998.
- [7] K. Chakrabarti and S. Mehrotra. The hybrid tree: An index structure for high dimensional feature spaces. *Proceedings of the IEEE International Conference on Data Engineering*, March 1999.
- [8] K. Chakrabarti and S. Mehrotra. Local dimensionality reduction: A new approach to indexing high dimensional spaces. *Technical Report, TR-MARS-00-04, University of California at Irvine*, 2000.
- [9] K. Chakrabarti, K. Porkaew, and S. Mehrotra. Supporting query refinement in multimedia databases. *Technical Report, MARS-TR-99-06*, 1999.
- [10] K. Chakrabarti, K. Porkaew, and S. Mehrotra. Efficient query refinement for multimedia similarity retrieval. *Proc. of ICDE*, 2000.
- [11] P. Ciaccia, M. Patella, and P. Zezula. M-tree: An efficient access method for similarity search in metric spaces. *Proc. of VLDB*, 1997.
- [12] R. Duda and P. Hart. Pattern classification and scene analysis. *Wiley, New York*, 1973.
- [13] M. Ester, J. Kohlhammer, and H. Kriegel. The dc-tree: A fully dynamic index structure for data warehouses. *Proc. of ICDE*, 2000.
- [14] M. Ester, H. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. *Proc. of KDD Conference*, 1996.
- [15] R. Fagin. Fuzzy queries in multimedia database systems. *Proceedings of PODS*, 1998.
- [16] C. Faloutsos. Fast searching by content in multimedia databases. *Data Engineering Bulletin* 18(4), 1995.
- [17] C. Faloutsos, W. Equitz, M. Flickner, W. Niblack, D. Petkovic, and R. Barber. Efficient and effective querying by image content. In *Journal of Intelligent Information Systems*, Vol. 3, No. 3/4, pages 231–262, July 1994.
- [18] C. Faloutsos and K.-I. D. Lin. Fastmap: A fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. In *Proc. ACM SIGMOD*, pages 163–174, May 1995.
- [19] U. Fayyad, C. Reina, and P. Bradley. Initialization of iterative refinement clustering algorithms. *Proc. of KDD*, 1998.
- [20] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Academic Press, second edition edition, 1990.
- [21] V. Ganti, R. Ramakrishnan, J. Gehrke, A. Powell, and J. French. Clustering large datasets in arbitrary metric spaces. *Proc. of ICDE*, 1999.
- [22] T. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 1985.
- [23] J. Gray and A. Reuter. *Transaction Processing: Concepts and Techniques*. Morgan Kaufmann, San Mateo, CA, 1993.
- [24] S. Guha, R. Rastogi, and K. Shim. Cure: An efficient clustering algorithm for large databases. *Proc. of SIGMOD*, 1998.
- [25] A. Guttman. R-trees: A dynamic index structure for spatial searching. In *Proc. ACM SIGMOD Conf.*, pp. 47–57., 1984.
- [26] G. R. Hjaltason and H. Samet. Ranking in spatial databases. *Proceedings of SSD*, 1995.
- [27] K. V. R. Kanth, D. Agrawal, and A. K. Singh. Dimensionality reduction for similarity searching dynamic databases. *Proc. of SIGMOD*, 1998.
- [28] N. Katayama and S. Satoh. The sr-tree: An index structure for high dimensional nearest neighbor queries. *Proc. of SIGMOD*, 1997.
- [29] L. Kaufman and P. Rousseeuw. Finding groups in data: An introduction to cluster analysis. *John Wiley and Sons*, 1990.

- [30] F. Korn, N. Sidiropoulos, and C. Faloutsos. Fast nearest neighbor search in medical image databases. *Proc. of VLDB*, 1996.
- [31] K. Lin, H. V. Jagadish, and C. Faloutsos. The TV-tree - an index structure for high dimensional data. In *VLDB Journal*, 1994.
- [32] N. Megiddo and A. Tamir. On the complexity of locating linear facilities in the plane. *Operation Research Letters*, 1982.
- [33] T. Mitchell. Machine learning. *McGraw Hill*, 1997.
- [34] R. Motwani and P. Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- [35] R. Ng and J. Han. Efficient and effective clustering methods for spatial data mining. *Proc. of VLDB*, 1994.
- [36] R. Ng and A. Sedighian. Evaluating multidimensional indexing structures for images transformed by principal component analysis. *Proc. of SPIE Conference*, 1996.
- [37] M. Ortega, Y. Rui, K. Chakrabarti, S. Mehrotra, and T. Huang. Supporting similarity queries in mars. *Proc. of ACM Multimedia 1997*, 1997.
- [38] K. Porkaew, K. Chakrabarti, and S. Mehrotra. Query refinement for content-based multimedia retrieval in MARS. *Proceedings of ACM Multimedia Conference*, 1999.
- [39] N. Roussopoulos, S. Kelley, and F. Vincent. Nearest neighbor queries. *Proceedings of SIGMOD*, 1995.
- [40] N. Roussopoulos, Y. Kotidis, and M. Roussopoulos. Cubetree: Organization of and bulk incremental updates on the data cube. *Proc. of SIGMOD*, 1997.
- [41] A. Szalay, P. Kunszt, A. Thakar, and J. Gray. Designing and mining multi-terabyte astronomy archives: The sloan digital sky survey. *Proc. of SIGMOD*, 2000.
- [42] M. Thomas, C. Carson, and J. Hellerstein. Creating a customized access method for blobworld. *Proc. of ICDE*, 2000.
- [43] R. Weber, H. Schek, and S. Blott. A quantitative analysis and performance study for similarity-search methods in high dimensional spaces. *Proc. of VLDB*, 1998.
- [44] D. White and R. Jain. Similarity indexing with the ss-tree. *Proc. of ICDE*, 1995.
- [45] T. Zhang, R. Ramakrishnan, and M. Livny. Birch: An efficient data clustering method for very large databases. *Proc. of SIGMOD*, 1996.

The material in this appendix can be read at the discretion of the reviewer and has been included only for the purpose of completeness.

A Principal Component Analysis

PCA examines the variance structure in the data and determines the directions along which the data exhibits high variance. The first principal component accounts for as much of the variability in the data as possible, and each succeeding component accounts for as much of the remaining variability as possible. Figure 15 shows a set of points and the two principal components. Since the first few principal components account for most of the variation in the data, the rest can be eliminated without significant loss of information. For example, in Figure 15, the second principal component can be eliminated, thus reducing the dimensionality from 2 to 1. The 1-d images of the 2-d points are obtained by projecting them on the first principal component (shown by squares in Figure 15). The reduced dimensional points are then indexed using an index structure.

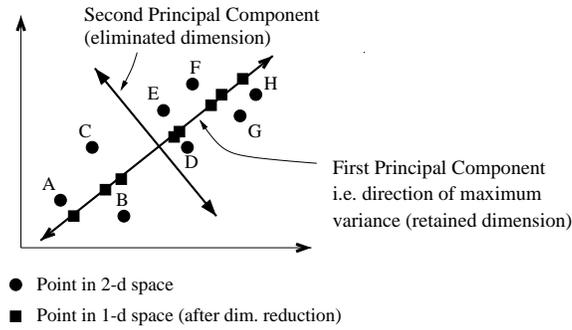


Figure 15: Global Dimensional Reduction (PCA or K-L Transform technique) where $D=2$, $d=1$.

We now describe how the principal components are computed algebraically. Let A be the $N \times D$ data matrix whose each row corresponds to a point in the original D -dimensional space. The first principal component is the eigenvector corresponding to the largest eigenvalue of the variance-covariance matrix of A , the second component correspond to the eigenvector with the second largest eigenvalue and so on. The mapping (to reduced dimensionality) corresponds to the well known Singular Value Decomposition (SVD) of data matrix A and can be done in $O(ND^2)$ time.

B Lower Bounding Lemma

Lemma 3 (Lower Bounding Lemma) $\mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i))$ always lower bounds $\mathcal{D}(P, Q)$.

Proof: Let P_i denote $\text{Image}(P, S_i)$ and Q_i denote $\text{Image}(Q, S_i)$. Let $P' = \bar{\Sigma}_{j=1}^D (P \bullet \Phi_i^{(j)})$ and $Q' = \bar{\Sigma}_{j=1}^D (Q \bullet \Phi_i^{(j)})$. Then, $\mathcal{D}(P', Q') = \mathcal{D}(P, Q)$ since Φ_i is orthonormal. Now,

$$P' = P_i + \overline{\text{ReconVect}}(P, S_i) + \bar{\Sigma}_{j=d_i+1}^D C_i^{(j)} \Phi_i^{(j)} \quad (10)$$

$$Q' = Q_i + \overline{\text{ReconVect}}(Q, S_i) + \bar{\Sigma}_{j=d_i+1}^D C_i^{(j)} \Phi_i^{(j)} \quad (11)$$

The vector distance $\overline{\text{Dist}}(P', Q')$ between P' and Q' is

$$\overline{\text{Dist}}(P', Q') = \overline{\text{Dist}}(P_i, Q_i) + (\overline{\text{ReconVect}}(P, S_i) - \overline{\text{ReconVect}}(Q, S_i)) \quad (12)$$

$$\Rightarrow \mathcal{D}(P', Q') = [\mathcal{D}(P_i, Q_i)^p + \|\overline{\text{ReconVect}}(P, S_i) - \overline{\text{ReconVect}}(Q, S_i)\|^p]^{1/p} \quad (13)$$

Since L_p functions obey triangle inequality,

$$\| \overline{\text{ReconVect}}(P, S_i) - \overline{\text{ReconVect}}(Q, S_i) \|_p \geq |(ReconDist(P, S_i, \mathcal{D}) - ReconDist(Q, S_i, \mathcal{D}))| \quad (14)$$

$$\Rightarrow \mathcal{D}(P', Q') \geq [\mathcal{D}(P_i, Q_i)^p + |(ReconDist(P, S_i, \mathcal{D}) - ReconDist(Q, S_i, \mathcal{D}))|^p]^{1/p} \quad (15)$$

Now,

$$\mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i)) = [\mathcal{D}(P_i, Q_i)^p + |(ReconDist(P, S_i, \mathcal{D}) - ReconDist(Q, S_i, \mathcal{D}))|^p]^{1/p} \quad (16)$$

Since $\mathcal{D}(P', Q') = \mathcal{D}(P, Q)$ and from Equations 15 and 16,

$$\mathcal{D}(Q, P) \geq \mathcal{D}(\text{NewImage}(P, S_i), \text{NewImage}(Q, S_i)) \quad (17)$$

■

C Range Query Algorithm

RangeSearch(Query $Q = \langle Q, \rho, \mathcal{D} \rangle$)	
1	for (i=1; i ≤ K; i++)
2	$Q_i \leftarrow \text{NewImage}(Q, S_i)$;
3	$Q_i \leftarrow \langle Q_i, \rho, \mathcal{D} \rangle$;
4	RangeSearchOnClusterIndex ($Q_i, R_i, S_i, result$);
5	for each $O \in \mathcal{O}$
6	if $\mathcal{D}(Q, O) \leq \rho$ $result \leftarrow result \cup O$;
RangeSearchOnClusterIndex(Query Q, Node T, Cluster S, Set $result$)	
1	if (T is a non-leaf node)
2	foreach child N of T
3	if $MINDIST(Q, N, \mathcal{D}) \leq \rho$ RangeSearchOnClusterIndex ($Q, N, S, result$);
4	else /* T is a leaf node */
5	for each object O in T
6	if $\mathcal{D}(Q, O) \leq \rho$
7	if $\mathcal{D}(\text{Original}(Q, S), \text{Original}(O, S)) \leq \rho$ $result \leftarrow result \cup O$;

Table 5: Range Query.

D Optimality of k nearest neighbor algorithm

Lemma 4 (Optimality of k-NN algorithm) *The k-NN algorithm is optimal.*

Proof: Let $\alpha = \max_{O \in \mathcal{A}} \mathcal{D}(Q, O)$ where \mathcal{A} is the set of final answers (the k nearest neighbors). The algorithm is optimal if it does not explore any indexed object O (in any cluster) (13-15) such that $\mathcal{D}(\text{NewImage}(O, S),$

$NewImage(Q, S) > \alpha$. Let us assume that it does explore such an object O . When O is explored, $|result| < k$ because otherwise the algorithm would have terminated before reaching this point. We will show that when O is explored, $|result|$ is at least k and hence prove the lemma (by contradiction). Each $O' \in \mathcal{A}$ has been explored before O since $\mathcal{D}(NewImage(O', S), NewImage(Q, S)) \leq \alpha < \mathcal{D}(NewImage(O, S), NewImage(Q, S))$ (by property of min-priority queue). Now $top.dist = \mathcal{D}(NewImage(O, S), NewImage(Q, S))$ when O is explored i.e. $top.dist > \alpha$. Since each $O' \in \mathcal{A}$ satisfies the condition $\mathcal{D}(Q, O) \leq \alpha$, it satisfies the condition $\mathcal{D}(Q, O) < top.dist$ and is hence added to $result$ (Line 7). So $|result|$ is at least k . ■

E Analysis of the height and balance of the global index structure

Let h_{GI} denote the the height of GI. Let h_{orig} denote the height of the original space index i.e. index on the entire dataset in the D -d original space. We assume that the multidimensional index structure used as the original space index is same as the one used to index the clusters (e.g., hybrid tree in both cases). Then, $h_{GI} \leq 1 + h_{orig}$. Since I_i is built on a subset of points of the entire set (i.e. $n_i \leq N$) and fewer dimensions (i.e. $d_i \leq D$), its height h_{I_i} cannot be greater h_{orig} . Since $h_{GI} = 1 + \max_i h_{I_i}$ and $h_{I_i} \leq h_{orig}$ for all i , $h_{GI} \leq 1 + h_{orig}$. The bound is a conservative one as the h_{GI} is usually smaller than h_{orig} due to the reduced size of the index.

We now show that GI is almost height-balanced. There are two factors that affect the height of a cluster index I_i : the number of points n_i and the subspace dimensionality d_i . Lower the value of n_i , lower the height. Also, lower the value of d_i , lower the height. Let I_{short} be the shortest index. Note $n_{short} \geq MinSize$. Let C_{short} and F_{short} denote the average number of entries in a leaf and index node of I_{short} respectively. Then, as explained in [23], the minimum possible height of I_{short} is $(1 + \lceil \log_{F_{short}}(\lceil \frac{MinSize}{C_{short}} \rceil) \rceil)$. Similarly, the maximum possible height of tallest index I_{tall} is $(1 + \lceil \log_{F_{tall}}(\lceil \frac{N}{C_{tall}} \rceil) \rceil)$ since $n_{tall} \leq N$. For space partitioning index structures (which is preferred for high dimensional indexing due to its ‘‘dimensionality-independent’’ fanout), $F_{short} \sim F_{tall}$ (say, F) [7]. C_{short} and C_{tall} depend on the respective subspace dimensionalities i.e. $\frac{C_{short}}{C_{tall}} \sim \frac{d_{tall}}{d_{short}}$. The maximum difference l_{max} in the lengths of *any* two paths from R to leaf is $l_{max} \sim \log_F(\frac{N * C_{short}}{MinSize * C_{tall}})$ i.e. $l_{max} \sim \log_F(\frac{N * d_{tall}}{MinSize * d_{short}})$. Usually, the subspace dimensionalities are close i.e. $d_{tall} \sim d_{short}$. For space-partitioning indexes, F is typically around 50-100 [7]. Under the above assumptions, $l_{max} \leq 1$ if $MinSize \geq \frac{N}{50}$ and $l_{max} \leq 2$ if $MinSize \leq \frac{N}{2500}$. In other words, with a proper size bound, l_{max} is usually 1 or at most 2, implying that GI is almost height balanced.

F Synthetic Data Generation

In order to generate the synthetic data, we use a method similar to that discussed in [45] but appropriately modified so that we can generate the different clusters in subspaces of different orientations and dimensionalities. The input parameters to the data generator is shown in Table 6. The generator generates k clusters with a total of $n \cdot (1 - o)$ points distributed among them using a Zipfian distribution with value z_{size} . The subspace dimensionality of each cluster also follows a Zipfian distribution with value z_{dim} , the average subspace dimensionality being d . Each cluster is generated as follows. For a cluster with size n_i and subspace dimensionality d_i (computed using the Zipfian distributions described above), we randomly choose d_i dimensions among the D dimensions as the subspace dimensions and generate n_i points in that d_i -d plane. Along each of the remaining $(D - d_i)$ non-subspace dimensions, we assign a randomly chosen coordinate to all the n_i points in the cluster. Let f_j be the randomly chosen coordinate along the j th non-subspace dimension. In the subspace, the points are spatially clustered into several regions (c regions on average) with each region having a randomly chosen centroid and an extent of r from the centroid along each of the d_i dimensions. After all the points in the cluster are generated, each

Parameter	Description	Default Value
n	Total number of points	100000
D	Original Dimensionality	64
k	Number of clusters	5
d	Average subspace dimensionality	10
z_{dim}	Skew in subspace dimensionality across clusters	0.5
z_{size}	Skew in size across clusters	0.5
c	Number of spatial clusters per cluster	10
r	Extent of a spatial cluster from centroid along each subspace dimension	0.5
p	Maximum displacement of points along each non-subspace dimension	0.1
o	Fraction outliers	0.05

Table 6: Input parameters to Synthetic Data Generator

point is displaced by a distance of at most p in either direction along each non-subspace dimension i.e. the point is randomly placed somewhere between $(f_j - p)$ and $(f_j + p)$ along the j th non-subspace dimension. The amount of displacement (i.e. value of p) determines the degree of correlation (since r is fixed). Lower the value, more the correlation. To make the subspaces arbitrarily oriented, we generate a random orthonormal rotation matrix (generated using MATLAB) and rotate the cluster by multiplying the data matrix with the rotation matrix. After all the clusters are generated, we randomly generate $N.o$ points (with random values along all D dimensions) as the outliers. The default values of the various parameters is shown in Table 6.

G Sensitivity to *MaxReconDist* parameter

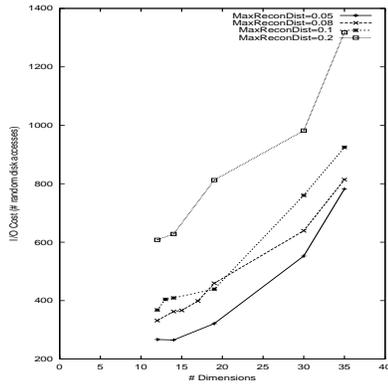


Figure 16: Sensitivity of I/O cost of LDR technique to *MaxReconDist*.

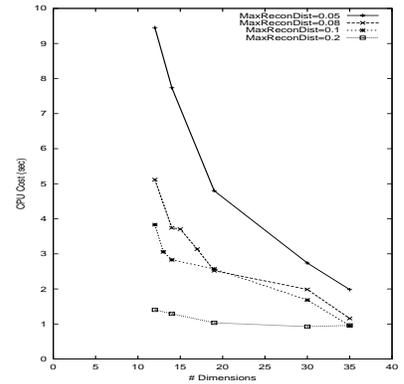


Figure 17: Sensitivity of CPU cost of LDR technique to *MaxReconDist*.

Figures 16 and 17 shows the sensitivity of the LDR technique to the *MaxReconDist* parameter in terms of I/O and CPU costs respectively. The I/O cost improves with decrease in *MaxReconDist* due to decrease in the information loss (i.e. fewer false positives) and hence decrease in post processing cost. However, with the decrease in *MaxReconDist*, the number of outliers increase as fewer points satisfy the reconstruction distance bound which causes the CPU cost to increase (the cost of scanning the outlier set) as shown in the Figure 17. The choice of *MaxReconDist* must consider the combined I/O and CPU cost; for example, $MaxReconDist = 0.08$ represents a good choice for this real-life dataset.