

# Fast feature selection using fractal dimension

*Caetano Traina Jr.*<sup>1</sup>    *Agma Traina*<sup>1</sup>    *Leejay Wu*<sup>2</sup>    *Christos Faloutsos*<sup>2</sup>

<sup>1</sup> Department of Computer Science and Statistics - University of São Paulo at São Carlos - Brazil

<sup>2</sup> Department of Computer Science - Carnegie Mellon University - USA  
{agma | caetano | lw2j | christos}@cs.cmu.edu

## Abstract

Dimensionality curse and dimensionality reduction are two issues that have retained high interest for data mining, machine learning, multimedia indexing, and clustering. We present a fast, scalable algorithm to quickly select the most important attributes (dimensions) for a given set of  $n$ -dimensional vectors. In contrast to older methods, our method has the following desirable properties: (a) it does not do rotation of attributes, thus leading to easy interpretation of the resulting attributes; (b) it can spot attributes that have nonlinear correlations; (c) it requires a constant number of passes over the dataset; (d) it gives a good estimate on how many attributes we should keep.

The idea is to use the ‘fractal’ dimension of a dataset as a good approximation of its intrinsic dimension, and to drop attributes that do not affect it. We applied our method on real and synthetic datasets, where it gave fast and good results.

## 1 - Introduction and Motivation

When managing the increasing volume of data which is generated by the organizations, a question which frequently arises is: “what part of this data is really relevant to be kept?”. Notice that usually the relations of the database have many attributes which are correlated with the others.

Attribute selection is a classic goal, as well as battling the “dimensionality curse” [Berchtold\_1998] [Pagel\_2000]. A careful chosen subset of attributes improves the performance and efficacy of a variety of algorithms. This is particularly true with redundant data, as many datasets can largely be well-approximated in fewer dimensions. This can also be seen as a way to compress data, as only the attributes which maintain the essential characteristics of the data are kept [Fayyad\_1998].

In this paper we introduce a novel technique that can discover how many attributes are significant to characterize a dataset. We also present a fast, scalable algorithm to quickly select the most significant attributes of a dataset. In contrast to other methods, such as Singular Value Decomposition (SVD) [Faloutsos\_1996], our method has the following desirable properties:

- (a) it does not rotate attributes, leading to easy interpretation of the resulting attributes;
- (b) it can spot attributes that have nonlinear and even non-polynomial correlations;
- (c) it is linear on the number of objects in the dataset;

---

<sup>1</sup>On leave at Carnegie Mellon University. This research has been funded by FAPESP (São Paulo State Foundation for Research Support - Brazil, under Grants 98/05556/5 and 98/0559-7).

<sup>2</sup>This material is based upon work supported by the National Science Foundation under Grants No. IRI-9625428, DMS-9873442, IIS-9817496, and IIS-9910606. Additional funding was provided by donations from NEC and Intel. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

(d) it gives a good estimate on how many attributes we should keep.

The main idea is to use the ‘fractal’ dimension of the dataset, and to drop attributes which do not affect it. The fractal dimension ( $D$ ) is relatively unaffected by redundant attributes, and our algorithm can compute it in **linear** time with respect to the number of objects. Thus, we propose a kind of backward-elimination algorithm to take advantage of the fast  $D$  computation. This algorithm sequentially removes attributes which contribute minimally to  $D$ .

The remainder of the paper is structured as follows. In the next section, we present a brief survey on the related techniques. Section 3 introduces the concepts needed to understand the proposed method. Section 4 presents the fractal dimension algorithm developed as well as the datasets used in the experiments. Section 5 gives the proposed method for attribute selection. Section 6 discusses the experiments and evaluation of the proposed method. Section 7 gives the conclusions of this paper.

## 2 - Survey

Numerous selection methods have been studied, including genetic algorithms; sequential feature selection algorithms such as forwards, backwards and bidirectional sequential searches; and feature weighting [Aha\_1995] [Scherf\_1997] [Vafaie\_1993]. A recent survey on attribute selection using machine learning techniques is presented in [Blum\_1997].

The singular value decomposition (SVD) technique provides another way of reducing the dimensionality of data by generating an ordered set of additional axes [Faloutsos\_1996]. However, this is not attribute selection, but instead axis generation as SVD returns vectors that do not need to correspond to the original attributes. These vectors may be inappropriate for assorted situations, such as those involving the presentation of data for human understanding; tasks where accessing additional attributes may be expensive; and when creating a training set to derive a classifier.

A common research challenge in attribute selection methods so far is the exponential growth of computing time required [Blum\_1997]. Indeed the induction methods proposed so far had super-linear or exponential computational complexity [Langley\_1997], as is the case with nearest neighbors, learning decision trees [John\_1994] [Kira\_1992], and Bayesian Networks [Singh\_1995]. Notice that these approaches are highly sensitive to both the number of irrelevant or redundant features present in the dataset, and to the size of the dataset, avoiding the use of samples [Langley\_1997].

Fractal dimension has been a useful tool for the analysis of spatial access methods [Belussi\_1995] [Kamel\_1994], indexing [Böhm\_2000], join selectivity estimation [Faloutsos\_2000], and analysis of metric trees [Traina\_2000]. However, to the best of the knowledge of the authors, it was never used to attribute selection.

### 3 - Fundamental Concepts

The most common way to store data is through tables with as many columns as there are features represented in the data, and as many lines as there are data elements. In this paper we are calling these tables as datasets, the features as attributes, and the data elements (or objects) as points in the space of features. In this way, a dataset is seen as points in an  $E$ -dimensional space, where  $E$  is the number of attributes.

We are especially interested in datasets which describe complex data, usually composed of numerical attributes. Features extracted

from images are well-known examples of high-dimensional datasets which are used in content-based image retrieval systems. For these datasets, it is difficult to choose the set of attributes that can be assigned as keys of the dataset. In this way, if one is interested in creating an index structure for the dataset the whole set of attributes needs to be considered.

Symbols	Definitions
$E$	embedding dimension (Euclidean dimensionality)
$D$	fractal dimension (intrinsic dimensionality)
$N$	number of points in the dataset
$C_{r,i}$	count ('occupancies') of points in the $i$ -th grid cell of side $r$
$r$	side of the grid cell
$S(r)$	total of occupancies for a specific grid cell side $r$
$R$	number of sides $r$ to plot $S(r)$

This leads to the previously mentioned dimensionality curse. Table 1- Definition of symbols

#### 3.1 - 'Embedding' and 'intrinsic dimensionality'

Our objective in this paper is to find a subset of the attributes that can be discarded when creating indexes or applying data mining techniques over the data, without compromising the results. Attributes that can be calculated from others are immediate candidates to discard, if the way to calculate them is known. However, in general, this correlation is not known. Thus, our objective turns into detection of correlations between attributes in a dataset, and how many redundant attributes the dataset has. This leads to the definition of the embedding and intrinsic dimensions.

**Definition 1** - The embedding dimension  $E$  of a dataset is the dimension of its address space. In other words, it is the number of attributes of the dataset.

The dataset can represent a spatial object that has a dimension lower than the space where it is embedded. For example, a line has an intrinsic dimensionality one, regardless if it is in a higher dimensional space.

**Definition 2** - The intrinsic dimension of a dataset is the dimension of the spatial object represented by

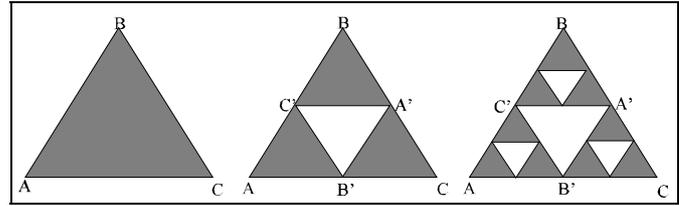
the dataset, regardless of the space where it is embedded.

Conceptually, if a dataset has all of its variables independent from the others, then its intrinsic dimension is the embedding dimension. However, whenever there is a correlation between two or more variables, the intrinsic dimensionality of the dataset is reduced accordingly. For example, each polynomial correlation (linear, quadratic, etc. ) reduces the intrinsic dimension by a unit. Other types of correlations can reduce the intrinsic dimension by different amount, even by fractional amounts, as we will show later.

Usually the embedding dimensionality of the dataset hides the actual characteristics of the dataset, and in general correlations between the variables in real datasets are not known and even the existence of correlations is not known either. This motivated us to look for a technique that allows one to find the intrinsic dimension of the dataset even when the existence of correlations is not identified. Knowing its intrinsic dimension, it is possible to decide how many attributes are in fact required to characterize a dataset.

### 3.2 - Fractals and Fractal Dimension

A fractal dataset is known by its characteristic of being self-similar. This means that the dataset has roughly the same properties for a wide variation in scale or size, i.e., parts of any size of the fractal are similar (exactly or statistically)



**Figure 1** - Recursive construction of the Sierpinsky triangle.

to the whole fractal. This idea is illustrated in Figure 1, which shows the first three steps to build the Sierpinsky triangle, a well-known point-set fractal. The Sierpinsky triangle is constructed from an equilateral triangle ABC, excluding its middle triangle A'B'C' and recursively repeating this procedure for each of the resulting smaller triangles. The Sierpinsky triangle is generated after infinite iterations of this procedure. The Sierpinsky triangle has an infinite perimeter, so it is not a 1-dimensional object. And it has no area, so it is not a 2-dimensional object. In fact, it has an intrinsic dimension, equal to  $\log 3 / \log 2 = 1.58$  [Schroeder\_1991]. For a real set of points, we measure the fractal dimension with the box-count plot, which is the basis of the algorithm to be proposed in Section 4.

**Definition 3 (Correlation Fractal dimension):** Given a dataset that has the self-similarity property in the range of scales  $\{r_1, r_2\}$ , its Correlation Fractal dimension  $D_2$  for this range is measured as

$$D_2 \equiv \frac{\partial \log \sum_i C_{r,i}^2}{\partial \log r}, \quad r \in [r_1, r_2]$$

From now on, we will use the correlation fractal dimension  $D_2$  as the intrinsic dimension  $D$ .

**Observation 1** - *For Euclidean objects, their fractal dimension corresponds to their Euclidean dimension, and the fractal dimension of Euclidean objects is always an integer number.*

For example, lines, circumferences and all standard curves have  $D=1$ ; planes, circles, squares and surfaces have  $D=2$ ; Euclidian volumes have  $D=3$ , and so on. Indeed, a line segment in any  $n$ -dimensional space will always have  $D=1$ , as well as a square will always have  $D=2$  even if the points are in a higher-dimensional space.

**Observation 2** - *The fractal dimension of a dataset cannot be greater than its embedding dimension.*

Many real datasets are fractals [Traina\_1999] [Schroeder\_1991]. Thus, for these datasets we can take the advantage of working with their correlation fractal dimension as their intrinsic dimension  $D$ . Table 1 summarizes the symbols used in this paper.

## 4 - Fractal Dimension Algorithm

This section presents an algorithm to compute the fractal dimension  $D$  of any given set of points in any  $E$ -dimensional space. A practical way to estimate  $D$  of a spatial dataset is using the box-counting approach [Schroeder\_1991]. Theoretically, this method gives a close approximation of the fractal dimension, and our experiments showed that it indeed does [Traina\_2000] [Traina\_1999]. One of the best published algorithm to calculate  $D$  of a dataset is an  $O(N \cdot \log(N))$  algorithm, where  $N$  is the number of points in the dataset [Belussi\_1995]. However, we developed a new, very fast,  $O(N)$  algorithm to implement it, which will be presented now.

Consider the address space of a point-set in an  $E$ -dimensional space, and impose an  $E$ -grid with grid-cells of side  $r$ . Focusing on the  $i$ -th cell, let  $C_{r,i}$  be the count ('occupancies') of points in each cell. Then, compute the value  $S(r) = \sum_i C_{r,i}^2$ . The fractal dimension is the derivative of  $\log(S(r))$  with

respect to the logarithm of the radius. As we assume self-similar datasets, we expect this derivative results in a constant value. Thus, we can obtain the fractal dimension  $D$  of a dataset plotting  $S(r)$  for different values of the radius  $r$ , and calculating the slope of the resulting line.

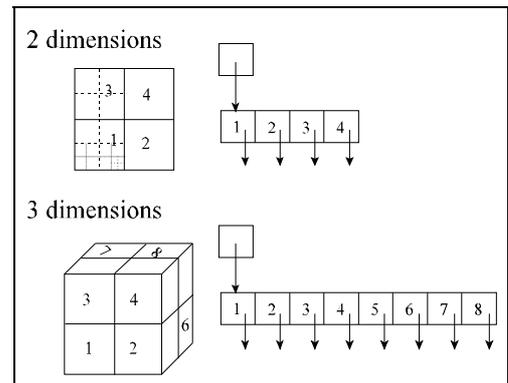
It is needed to process  $S(r)$  for many values of  $r$ . To avoid read the dataset again for each value of the radius, we propose to create a multi-level grid structure, where each level has a radius the half of the size of the previous level ( $r=1, 1/2, 1/4, 1/8$ , etc.). Each level of the structure corresponds to a different radius, so the depth of the structure is equal to the number of points in the resulting graph. This structure is created in main memory, so the number of points in the graph is limited by the amount of main memory available. If this graph is linear for a suitable range of radii, the dataset is a fractal and its

fractal dimension  $D$  is the slope of the fitting line of this graph.

The proposed algorithm is linear on the number of points in the dataset. The computational complexity of the algorithm is  $O(N * E * R)$ , where  $N$  is the number of objects in the dataset,  $E$  is the embedding dimensionality, and  $R$  is the number of points used to plot the  $S(r)$  function. This shows that the algorithm is scalable to datasets of any size.

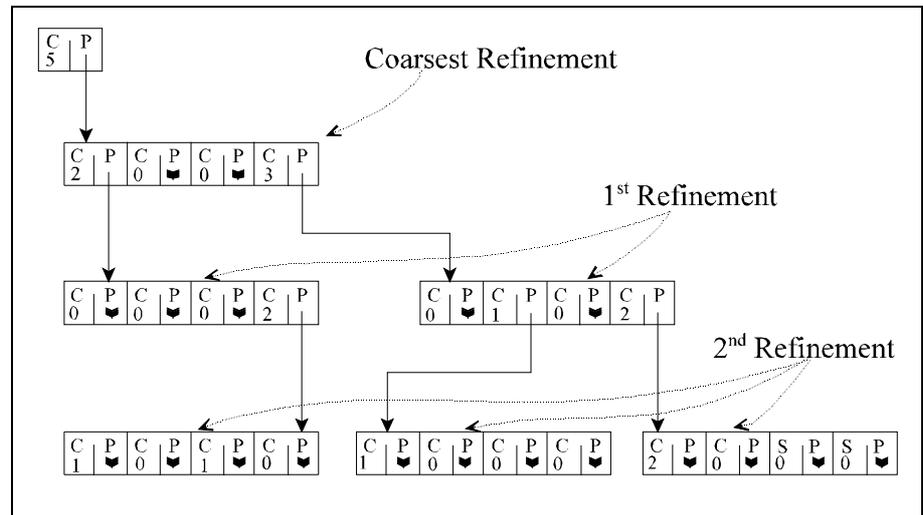
For each given resolution  $r$ , only the cells which have at least one already processed point are maintained, counting the sum of occupancies  $C_{r,i}$  of this cell. In this way, each new point is directly associated to a cell in each level, without the need to be compared with the previously read points. Figure 2 shows the structure used in the algorithm for 2- and 3-dimensional datasets.

The coarsest resolution of the space of points generates  $2^n$  cells. In the next level each cell is split into other  $2^n$  cells, and so on. Given that the position of each cell in the space is always known, each cell is represented by: the sum of occupancies  $C_{r,i}$  in this cell, and the pointers to the cells in the next level covered by this cell (see Figure 2). This structure is a kind of a multidimensional “quad-tree” (oct-tree,  $E$ -dim-tree). Figure 3 shows an example of this structure for a dataset with five points in three levels in a 2-dimensional space.



**Figure 2** - Representation of grid cells in 2- and 3-dimensional space.

Notice that new cells are added to the structure on demand. Thus, only cells occupied by at least one point are created ( $C_{r,i} > 0$ ). This algorithm processes the set of points only once, so it is indeed very fast. Figure 4 summarizes this algorithm.



**Figure 3** - Example of the data structure used for calculating the Sum of Occupancies of a dataset with 5 points (with three-level resolution).

As the grid resolution increases, the number of pointers to empty cells increases as well. Thus, for high-dimensional datasets it is worthwhile to keep the cells as linked lists instead of arrays. We implemented this structure as an object in C++, using an array for datasets with the embedding dimension less or equal three, and using a linked list for datasets with higher dimensionality.

**Algorithm 1** Compute fractal dimension  $D$  of a dataset  $A$  (box-count approach)  
input: normalized dataset  $A$  ( $N$  rows, with  $E$  dimensions/attributes each)  
output: fractal dimension  $D$

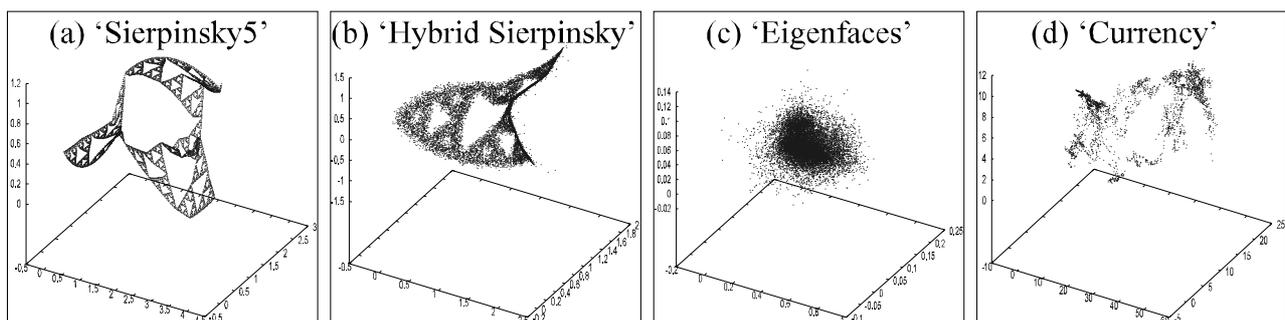
*Begin*  
For each desirable grid-size  $r=1/2^j, j= 1, 2, \dots, l$   
  For each point of the dataset  
    Decide which grid cell it falls in (say, the  $i$ -th cell)  
    Increment the count  $C_i$  ('occupancy')  
  Compute the sum of occupancies  
   $S(r) = \sum C_i^2$   
Print the values of  $\log(r)$  and  $\log(S(r))$  generating a plot;  
Return the slope of the linear part of the plot as the fractal dimension  $D$  of the dataset  $A$ .  
*End*

**Figure 4** - Correlation fractal dimension algorithm.

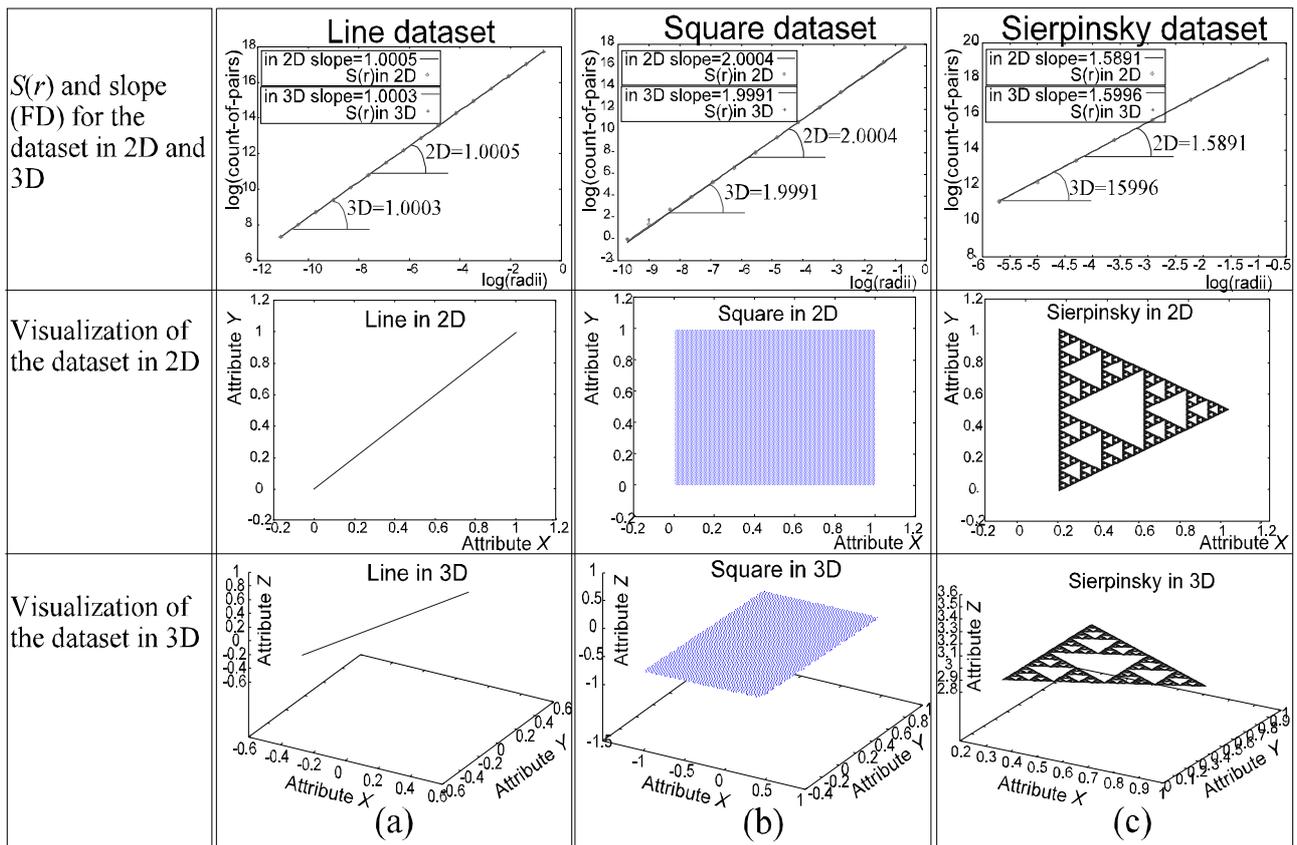
#### 4.1 - Datasets used in the experiments

We used synthetic and real datasets to evaluate our method. Figure 5 shows a mapping in a 3-dimensional space of the higher-dimensional datasets used in the experiments. This mapping was done through the *FastMap* algorithm [Faloutsos\_1995]. We used two synthetic datasets built over a Sierpinsky triangle (9,841 points in a 2D space), adding three more attributes to the dataset in order to test our method. The synthetic datasets are:

- “Sierpinsky5” (see Figure 5(a)) - The original 2D points of the original dataset ( $x, y$ ) became 5D points ( $a=x, b=y, c=a+b, d=a^2 + b^2, e=a^2 - b^2$ ). The three latest coordinates included in this dataset are strongly correlated with the two first coordinates. Thus, the fractal dimension (1.68) of the new dataset is close to the fractal dimension of the original Sierpinsky triangle.
- “Hybrid5” (see Figure 5(b)) - The original 2D points of the Sierpinsky triangle ( $x, y$ ) became 5D points ( $a=x, b=y, c=f(a,b), d=\text{random1}, e=\text{random2}$ ). As the two latest coordinates include random noise to the dataset the fractal dimension of this dataset is equal to 3.62, basically the dimensionality of the Sierpinsky (1.58) plus the dimensionality of a square in 2D (2.00). The third variable (' $c$ ') is non-linearly depending on the others. It is obtained by the algorithm during the Sierpinsky triangle generation.



**Figure 5** - Three-dimensional mappings of the datasets used in the experiments of the proposed method (“FDR”).



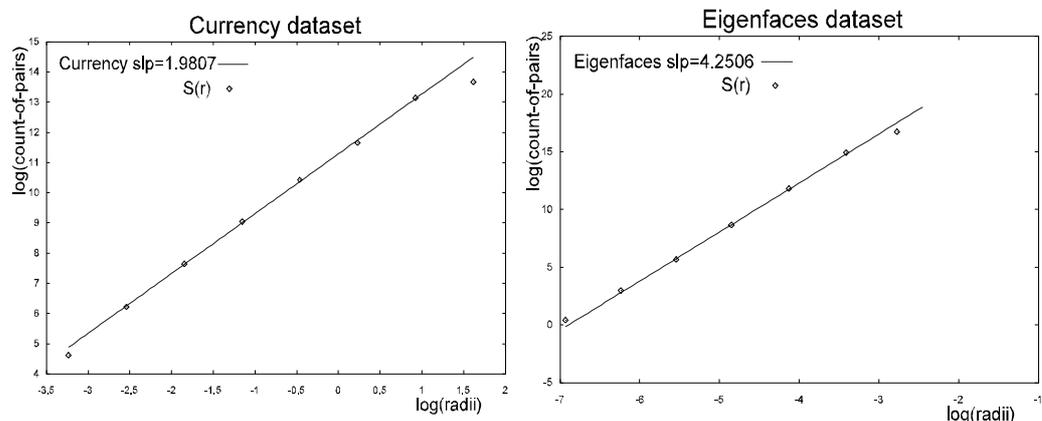
**Figure 6** - Fractal dimension of synthetic datasets embedded in 2- and 3-dimensional spaces. (a) Line; (b) Square; (c) Sierpinsky triangle.

Also two real datasets were used to evaluate our proposed method. Here are the datasets:

- “Currency” (see Figure 5(c)) - This is a 6-dimensional dataset, presenting the normalized exchange rate of currencies based on Canadian Dollar. The data was collected from 01/02/87 until 01/28/97. This resulted in  $N=2,561$  observations made on working days. Each attribute corresponds to a currency (a = Hong Kong Dollar, b = Japanese Yen, c = American Dollar, d = German Mark, e = French Franc, f = British Pound).

- “Eigenfaces” (see Figure 5(d)) - a dataset of 11,900 face vectors given by the Informat project

[Wactlar\_1996] at Carnegie Mellon University. Each face was processed



**Figure7** - Fractal dimension of real datasets. (a) ‘Currency’ dataset; (b) ‘Eigenfaces’ dataset.

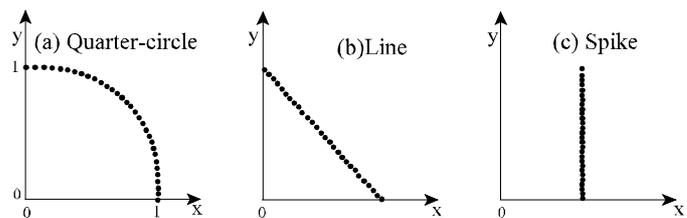
with the eigenfaces method [Turk\_1991], resulting in 16-dimensional vectors.

Now looking at Figure 6, we are able to show that Observation 1 stated previously indeed holds. It can be seen that the correlation fractal dimension proposed indeed gives the fractal dimension of the datasets regardless of their embedding dimension. Figure 7 shows the correlation fractal dimension of the real datasets used in this paper.

## 5 - Attribute Selection Algorithm

### 5.1 - Intuition

In this Section we present an approach to quickly discard some attributes (dimensions) from the original dataset, taking advantage of the fractal dimension concept. We stated in



**Figure 8** - Example of point sets in  $E=2$ -dimensional space.

Section 2 that the fractal dimension ( $D$ ) of a

(a) ‘Quarter-circle’, (b) ‘Line’, (c) ‘Spike’.

dataset cannot exceed the embedding dimensionality ( $E$ ) of the dataset. Moreover, there are  $\lceil D \rceil$  attributes which cannot be determined from the others. Since  $D \leq E$ , there are at least  $E - \lceil D \rceil$  attributes which can be correlated with the others. Correlated attributes contribute to increase the complexity of any treatment that the dataset must be submitted to, such as spatial indexing in a database, and knowledge retrieval in data mining processes. Moreover, the correlated attributes can be re-obtained from the other attributes. Hence, whenever it is possible, such attributes should be detected and dropped from the dataset.

**Definition 4 - Partial fractal dimension ( $pD$ )** : Given a dataset  $A$  with  $E$  attributes, this measurement is obtained through the calculation of the correlation fractal dimension of this dataset excluding one or more attributes from the dataset.

Figure 8(a) illustrates the intuition behind our approach. This is the ‘Quarter-circle’ dataset, which points are in  $E=2$  dimensions, and fractal dimension  $D=1$ . Notice that, the two attributes  $x$  and  $y$  are correlated in a nonlinear way  $y = \sqrt{1-x^2}$ . Also notice that the traditional dimensionality reduction method, SVD, only works well for linear correlations. Computing the fractal dimension  $D=1$ , gives a hint that maybe the two attributes are correlated. Thus, the points projected on one axes (say  $x$ ) probably will preserve the original distances. The fractal dimension of the projected points will reveal to us how well preserved the intrinsic properties of the dataset are. In this specific case, the  $pD$  for  $x$  is  $pD=0.9$  which means that the mode of the dataset was kept after projection. Consider also the Figure 8(b) and 8(c) presenting ‘Line’ and ‘Spike’ examples respectively. Again, our approach will correctly

flag attributes  $x, y$  for omission, but it will not allow to drop the attribute  $y$  in the ‘Spike’ picture.

## 5.2 - Proposed Algorithm - “FDR”

The algorithm to be described uses the approach of backward elimination of the attributes. We named it as *Fractal Dimension Reduction* (FDR). The proposed idea is to calculate the correlation fractal dimension of the whole dataset, and also to calculate its  $pD$  dropping one of its  $E$  attributes at a time. Thus, it will result in  $E$  partial fractal dimensions. The process continues selecting the attribute that leads to the minimum difference in the  $pD$  to the whole dataset. If this difference is within a small threshold, we can be confident that this attribute contributes almost nothing to the overall characteristics of the dataset. Therefore, this attribute can be dropped from the list of important attributes that characterizes the dataset. The threshold depends on how precise the resulting dataset needs to be to preserve the characteristics of the original dataset.

The algorithm is iterative, i.e., using the resulting set of attributes, repeat the previous reduction algorithm, until there are no more attributes to be dropped without changing the previous partial fractal dimension more than a fixed threshold.

If there are two or more attributes correlated, this algorithm will sequentially drop attributes using this correlation, until only the number of attributes that corresponds to independent attributes remain. For example, if there are three attributes  $\{a, b, c\}$ , where the third is a function of the previous two, e.g.,  $c=a+b$ , any of the three attributes can be dropped, because the others can be used to derive the one dropped. However, if there is no other correlation linking the remaining attributes, then no other attribute could be dropped without mischaracterizing the dataset. Figure 9 presents the algorithm used to generated the classification of the attributes which are presented ordered by their significance. That

### **Algorithm 2** - *Fractal dimensionality reduction (FDR) algorithm*

input: dataset  $A$

output: list of attributes in the reverse order of their importance

*Begin*

- 1- Compute the fractal dimension  $D$  of the whole dataset;
- 2- Initially set all attributes of the dataset as the significant ones, and the whole fractal dimension as the current  $D$ ;
- 3- While there are significant attributes do:
  - 4- For every significant attribute  $i$ , compute the partial fractal dimensions  $pD_i$  using all significant attributes excluding attribute  $i$ ;
  - 5- Sort the partial fractal dimensions  $pD_i$  obtained in step 4 and select the attribute  $a$  which leads to the minimum difference (current  $D - pD_i$ );
  - 6- Set the  $pD_i$  obtained removing attribute  $a$  as the current  $D$ ;
  - 7- Output attribute  $a$  and remove it from the set of important attributes;

*end*

**Figure 9** - Fractal dimensionality reduction (FDR) algorithm.

is, the first attribute to be dropped is the least important attribute, the second attribute dropped is the second most important attribute and so on.

## 6 - Experiments and Evaluation

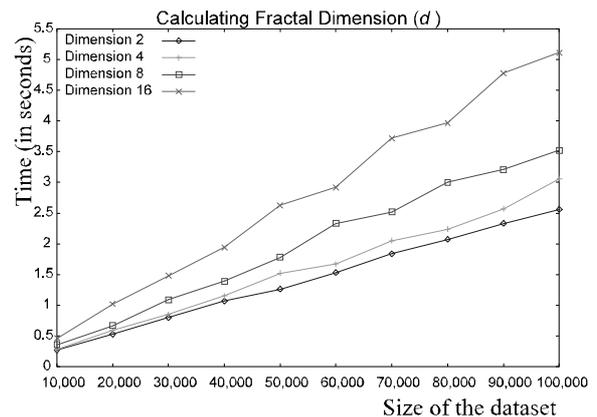
We did experiments to answer the following questions:

- (1) How scalable are the proposed algorithms?
- (2) How many attributes should be kept in order to reduce the dimensionality of the dataset?

The following sections will clarify these points. The experiments and measurements were taken on a 450 MHz Pentium II machine with 128 Mbytes of RAM under Windows NT4.0. All the proposed algorithms were implemented in C++ language.

### 6.1 - Scalability of the proposed method

The algorithm developed to obtain the correlation fractal dimension is linear over the number of points in the dataset, i.e.,  $O(N)$ . The embedding dimensionality  $E$  of the dataset is a constant involved in the process as well as the maximum resolution of the grid  $l$  (i.e., the number of grid sizes), then the complexity of our algorithm is  $O(N*R*E)$ . However,  $E$  and  $l$  are small values,  $R$  is typically equal 20 (value used in our experiments). These values are much smaller than the number of points in the datasets, which are in order of thousands. Figure



**Figure 10** - Wall-clock time (in seconds) needed to obtain the fractal dimension of varying sized datasets. The curves show the datasets with 2, 4, 8, and 16 dimensions.

10 shows the wall-clock time required to get the fractal dimension against the size of the dataset. The datasets have varying number of points in 2, 4, 8 and 16-dimensional spaces generating 20 grid sizes for each of which. Figure 10 shows that the execution time of this algorithm is linear on the number of points in the dataset.

The algorithm developed to select the attributes of a dataset by their significance is very fast. Instead of the super-linear time over the size of the

Dataset	Number of Points $N$	Embedding dimensionality $E$	Intrinsic Dimensionality $D$	time (in seconds)
'Sierpinsky5'	9,841	5	1.597	6.24
'Hybrid5'	9,841	5	3.627	7.03
'Eigenfaces'	11,900	16	4.250	132.82
'Currency'	2,561	6	1.980	2.54

**Table 2** - Wall-clock time (in seconds) spent to run the backward-selection algorithm on the datasets presented. A summarization of the datasets is also given.

dataset ( $N$ ) being analyzed, as it is needed by the machine learning techniques [Blum\_1997] our FDR algorithm is linear on  $N$  (number of objects) and quadratic on the embedding dimensionality ( $E$ ) of the dataset. Table 2 shows the wall-clock time needed to generate the classification of the attributes for the datasets we presented in this paper. Table 2 also summarizes the meaningful information of the datasets.

## 6.2 - Dimensionality reduction using fractal dimension

Figure 11 presents the graphs generated by the FDR algorithm on the test datasets. Figure 11(a) shows the graph of the  $pD$  of the ‘Sierpinsky5’ dataset when its attributes are sequentially dropped. From this plot, it can be seen that just two attributes are enough to characterize this dataset. Our algorithm drops  $c=a+b$ ,  $e=a^2-b^2$  and  $a$  attributes, holding  $b$  and  $d=a^2+b^2$ , with a resulting partial fractal dimension  $pD=1.568$  (versus a whole  $pD=1.597$ ). Notice that knowing  $b$  and  $d=a^2+b^2$ ,  $a$  and all other attributes can be recalculated.

Figure 11(b) presents the same plot of the  $pD$  for the ‘Hybrid5’ dataset when its attributes are sequentially dropped. Looking at this plot it can be seen that four attributes are needed to characterize this dataset. Just the  $c=f(a+b)$  attribute can be dropped, as every other attribute contributes with a significant portion of  $D$ . This is correct, as the attributes  $a$  and  $b$  correspond to the original Sierpinsky triangle points, and the attributes  $d$  and  $e$  depend on random numbers, which are independent variables and cannot be obtained from the other attributes. Also, as ‘Hybrid5’ dataset has  $D=3.62$  it is expected four attributes to remain.

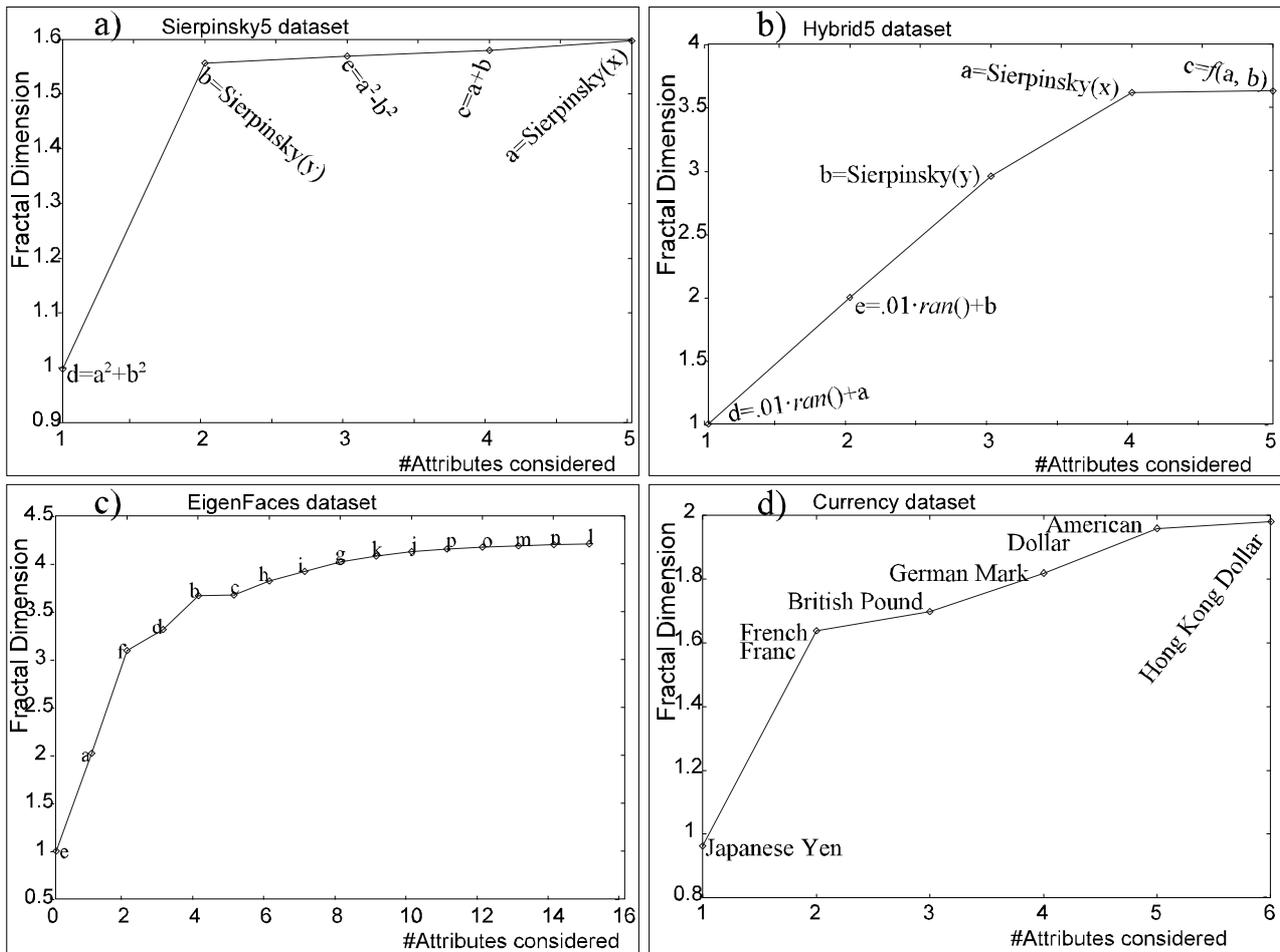
Figure 11(c) shows the plot of the  $pD$  of the ‘Eigenfaces’ dataset when its attributes are sequentially dropped. From this plot, we can see that, from the original 16 attributes, just five are enough to characterize this dataset  $\{b, d, f, a, e\}$ . The resulting partial fractal dimension with five attributes is 3.815, and the whole partial fractal dimension is 4.207 (that is, eleven attributes contribute only 0.392 to the whole fractal dimension).

Figure 11(d) shows the plot of the  $pD$  of the ‘Currency’ dataset when its attributes are sequentially dropped. It shows that the Hong Kong Dollar is the only currency that can be immediately dropped. This is correct, as we know that the Hong Kong Dollar is linked to the American Dollar, so there is some strong correlation between both currencies. The other currencies have more independent behaviors, as their contribution to the whole fractal dimension is a value between 0.16 and 0.68.

The following observations can now be made:

**Observation 3** - *The intrinsic dimensionality gives a lower bound of the number of attributes needed to keep the essential characteristics of the dataset.*

This means that at least the number of attributes equal to the ceiling function on  $D$  needs to remain.



**Figure 11** - Plots of the number of attributes dropped versus the partial fractal dimensions for the following datasets: (a) ‘Sierpinsky5’ (b) ‘Hybrid5’ (c) ‘Eigenfaces’ (d) ‘Currency’.

**Observation 4** - *The most independent attributes are saved to the end of the process.*

As this occurs (by the construction of the algorithm), the process can stop early, when the minimum number of attributes is achieved.

### 6.3 - Discussion

Intuitively the attribute selection could be performed in backward or forward direction. If there is only polynomial correlation between the attributes, both backward or forward selection works well. However, when there is a fractal correlation between the attributes (such as the  $x$  and  $y$  coordinates in the Sierpinsky triangle), the experiments showed that the backward selection works better.

The fractal dimension  $D$  is a guide to know when to stop the backward selection algorithm FDR. Indeed,  $\lceil D \rceil$  is the minimum number of attributes that must be in the resulting set. This is due to the fact these  $\lceil D \rceil$  attributes preserve the essential characteristics of the dataset.

## 7 - Conclusions

The main contribution of this paper is the proposal of a novel approach in feature selection and

dimensionality reduction, using the concept of fractal dimension. This approach leads to a method to reduce the dimensionality of spatial datasets and it has the following properties:

- it can detect the hidden correlations which exist in the dataset, spotting how many attributes strongly affect the behavior of the dataset regarding index and retrieval operations.
- it can show the attributes that have nonlinear and even non-polynomial correlations, where the traditional SVD method fails;
- it provides a small subset of attributes that can represent the whole dataset.
- it is scalable on the number  $N$  of elements in the dataset -  $O(N)$ . This is a striking advantage over methods from Machine Learning field [Blum\_1997] which are super-linear on the number of objects  $N$ .
- it can be applied to high-dimensional datasets as well.
- it does not rotate the address space of the dataset. Thus, it leads to easy interpretation of the resulting attributes.

Other contributions are:

- the detailed design of the single pass algorithm to compute the correlation fractal dimension of any spatial dataset. This algorithm is  $O(N)$  thus scaling up for arbitrary size datasets. This algorithm works in main memory, but the amount of memory available limits only the resolution of results, and not the size or dimension of the dataset.
- the quick backward attribute reduction algorithm. As it uses the quick algorithm to calculate the fractal dimension, it is also linear on the size of the dataset. Moreover, it can quickly compute the meaningful attributes (seconds), in contrast to current methods that take hours or days to give answers.
- experiments on synthetic and real datasets, showing the effectiveness and speed of the results.

## References

- [Aha\_1995] D. W. Aha and R. L. Bankert, "A Comparative Evaluation of Sequential Feature Selection Algorithms," in Proceedings of the Fifth Intl. Workshop on Artificial Intelligence and Statistics, Ft. Lauderdale, FL, 1995.
- [Belussi\_1995] A. Belussi and C. Faloutsos, "Estimating the Selectivity of Spatial Queries Using the 'Correlation' Fractal Dimension," in 21th Intl. Conf. on Very Large Data Bases (VLDB), Zurich, Switzerland, 1995.
- [Berchtold\_1998] S. Berchtold, C. Böhm, H.-P. Kriegel, "The Pyramid-Tree: Breaking the Curse of Dimensionality," in ACM SIGMOD Intl. Conf. on Management of Data, Seattle, Washington, 1998.
- [Blum\_1997] A. L. Blum and P. Langley, "Selection of Relevant Features and Examples in Machine Learning," *Artificial Intelligence*, vol. 97, pp. 245-271, 1997.
- [Böhm\_2000] C. Böhm and H.-P. Kriegel, "Dynamically Optimizing High-Dimensional Index Structures," in

7th Intl. Conf. on Extending Database Technology (EDBT), Konstanz, Germany, 2000.

- [Faloutsos\_1996] C. Faloutsos, *Searching Multimedia Databases by Content*. Boston: Kluwer Academic Publishers, 1996.
- [Faloutsos\_1995] C. Faloutsos and K.-I. Lin, “FastMap: A Fast Algorithm for Indexing, Data-Mining and Visualization of Traditional and Multimedia Datasets,” in ACM SIGMOD Intl. Conf. on Management of Data, San Jose, California, 1995.
- [Faloutsos\_2000] C. Faloutsos, B. Seeger, A. J. M. Traina, C. T. Jr., “Spatial Join Selectivity Using Power Laws,” in to be published in the ACM SIGMOD Intl. Conf. on Management of Data, Dallas, Texas, USA, 2000.
- [Fayyad\_1998] U. Fayyad, “Mining Databases: Towards Algorithms for Knowledge Discovery,” *Bulletin of the IEEE Technical committee on Data Engineering*, vol. 21, pp. 39-48, 1998.
- [John\_1994] G. H. John, R. Kohavi, K. Pfleger, “Irrelevant Features and the Subset Selection Problem,” in 11<sup>th</sup> Intl. Conf. on Machine Learning, New Brunswick, NJ, USA, 1994.
- [Kamel\_1994] I. Kamel and C. Faloutsos, “Hilbert R-tree: An Improved R-tree using Fractals,” in 20th Intl. Conf. on Very Large Data Bases (VLDB), Santiago de Chile, Chile, 1994.
- [Kira\_1992] K. Kira and L. L. Rendell, “A Practical Approach fo Feature Selection,” in 9<sup>th</sup> Intl. Conf. on Machine Learning, Aberdeen Scotland, 1992.
- [Langley\_1997] P. Langley and S. Sage, “Scaling to Domains with many Irrelevant Features,” in *Computational learning theory and natural learning systems*, vol. 4, R. Greiner, Ed. Cambridge, MA: MIT Press, 1997.
- [Pagel\_2000] B.-U. Pagel, F. Korn, C. Faloutsos, “Deflating the Dimensionality Curse Using Multiple Fractal Dimensions,” in 16th Intl. Conf. on Data Engineering (ICDE), San Diego, CA - USA, 2000.
- [Scherf\_1997] M. Scherf and W. Brauer, “Feature Selection by Means of a Feature Weighting Approach,” Technische Universität München, Munich 1997.
- [Schroeder\_1991] M. Schroeder, *Fractals, Chaos, Power Laws*, 6 ed. New York: W.H. Freeman and Company, 1991.
- [Singh\_1995] M. Singh and G. M. Provan, “A Comparison of Induction Algorithms for Selective and Non-selective Bayesian Classifiers,” in 12<sup>th</sup> Intl. Conf. on Machine Learning, Lake Tahoe, CA, USA, 1995.
- [Traina\_1999] C. Traina, A. J. M. Traina, C. Faloutsos, “Distance Exponent: A New Comcept for Selectivity Estimation in Metric Trees,” Carnegie Mellon University, Pittsburgh, PA CMU-CS-99-110, March 1, 1999 1999.
- [Traina\_2000] C. Traina, A. J. M. Traina, C. Faloutsos, “Distance Exponent: a New Concept for Selectivity Estimation in Metric Trees,” in IEEE Intl. Conf. on Data Engineering (ICDE), San Diego - CA, 2000.
- [Turk\_1991] M. Turk and A. Pentland, “Eigenfaces for Recognition,” *Journal of Cognitive Neuroscience*, vol. 3, pp. 71-86, 1991.
- [Vafaie\_1993] H. Vafaie and K. A. D. Jong, “Robust Feature Selection Algorithms,” in Intl. Conf. on Tools with AI, Boston, MA, 1993.
- [Wactlar\_1996] H. D. Wactlar, T. Kanade, M. A. Smith, S. M. Stevens, “Intelligent Access to Digital Video: Informedia Project,” *IEEE Computer*, vol. 29, pp. 46-52, 1996.