PART II

INDEXING SIGNALS

7 PROBLEM - INTUITION

7.1 INTRODUCTION

The problem we focus on is the design of fast searching methods that will search a database of multimedia objects, to locate objects that match a query object, exactly or approximately. Objects can be 2-dimensional color images, gray-scale medical images in 2-d or 3-d (eg., MRI brain scans), 1-dimensional time sequences, digitized voice or music, video clips etc. A typical query by content would be, eg., 'in a collection of color photographs, find ones with a same color distribution as a sunset photograph'.

Specific applications include image databases; financial, marketing and production time sequences; scientific databases with vector fields; audio and video databases, DNA/genome databases, etc. In such databases, typical queries would be 'find companies whose stock prices move similarly', or 'find images that look like a sunset', or 'find medical X-rays that contain something that has the texture of a tumor'.

Searching for similar patterns in such databases as the above is essential, because it helps in predictions, decision making, computer-aided medical diagnosis and teaching, hypothesis testing and, in general, in 'data mining' [AGI+92, AIS93b] [AIS93a, AS94, HS95] and rule discovery.

The first important step is to provide a measure for the distance between two objects. We rely on a domain expert to supply such a distance function $\mathcal{D}()$:

Definition 7.1 Given two objects, O_A and O_B , the distance (= dis-similarity) of the two objects is denoted by

$$\mathcal{D}(O_A, O_B) \tag{7.1}$$

For example, if the objects are two (equal-length) time sequences, the distance $\mathcal{D}()$ could be their Euclidean distance (sum of squared differences, see Eq. 7.2).

Similarity queries can been classified into two categories:

- Whole Match: Given a collection of N objects O_A, O_B, \ldots, O_N and a query object Q, we want to find those data objects that are within distance ϵ from Q. Notice that the query and the objects are of the same type: for example, if the objects are 512×512 gray-scale images, so is the query.
- **Sub-pattern Match:** Here the query is allowed to specify only part of the object. Specifically, given N data objects (eg., images) O_A, O_B, \ldots, O_N , a query (sub-)object Q and a tolerance ϵ , we want to identify the parts of the data objects that match the query. If the objects are, eg., 512×512 gray-scale images (like medical X-rays), in this case the query could be, eg., a 16×16 sub-pattern (eg., a typical X-ray of a tumor).

Additional types of queries include the 'nearest neighbors' queries (eg., 'find the 5 most similar stocks to IBM's stock') and the 'all pairs' queries or 'spatial joins' (eg., 'report all the pairs of stocks that are within distance ϵ from each other'). Both the above types of queries can be supported by our approach: As we shall see, we reduce the problem into searching for multi-dimensional points, which will be organized in R-trees; in this case, we know of algorithms for both nearest-neighbor search as well as spatial joins, as discussed in Chapter 5. Thus, we do not focus on nearest-neighbor and 'all-pairs' queries.

For all the above types of queries, the ideal method should fulfill the following requirements:

- It should be *fast*. Sequential scanning and distance calculation with each and every object will be too slow for large databases.
- It should be 'correct'. In other words, it should return all the qualifying objects, without missing any (i.e., no 'false dismissals'). Notice that 'false alarms' are acceptable, since they can be discarded easily through a post-processing step. Of course, as we see, eg., in Figure 8.2, we try to keep their number low, so that the total response time is minimized.

- The proposed method should require a small space overhead.
- The method should be dynamic. It should be easy to insert, delete and update objects.

As we see next, the heart of the proposed approach is to use k feature extraction functions, to map objects into points in k-dimensional space; thus, we can use highly fine-tuned database spatial access methods to accelerate the search. In the next Chapter we describe the details of the main idea. In Chapters 8 and 9 we describe how this idea has been applied for time sequences and color images. Chapter 10 discusses how to extend the ideas to handle sub-pattern matching in time sequences. Chapter 11 discusses a fast, approximate method of extracting features from objects, so that the distance is preserved. Chapter 12 lists the conclusions for this Part.

7.2 BASIC IDEA

To illustrate the basic idea, we shall focus on 'whole match' queries. There, the problem is defined as follows:

- we have a collection of N objects: O_A, O_B, \ldots, O_N
- the distance/dis-similarity between two objects (O_i, O_j) is given by the function $\mathcal{D}(O_i, O_j)$, which can be implemented as a (possibly, slow) program
- the user specifies a query object Q, and a tolerance ϵ

Our goal is to find the objects in the collection that are within distance ϵ from the query object. An obvious solution is to apply sequential scanning: For each and every object O_i $(1 \le i \le N)$, we can compute its distance from Q and report the objects with distance $\mathcal{D}(Q, O_i) \le \epsilon$.

However, sequential scanning may be slow, for two reasons:

1. the distance computation might be expensive. For example, the editing distance [HD80] in DNA strings requires a dynamic-programming algorithm, which grows like the product of the string lengths (typically, in the hundreds or thousands, for DNA databases).

2. the database size N might be huge.

Thus, we are looking for a faster alternative. The proposed approach is based on two ideas, each of which tries to avoid each of the two disadvantages of sequential scanning:

- a 'quick-and-dirty' test, to discard quickly the vast majority of non-qualifying objects (possibly, allowing some false-alarms)
- the use of Spatial Access Methods (SAMs), to achieve faster-than-sequential searching, as suggested by Jagadish [Jag91].

The case is best illustrated with an example. Consider a database of time sequences, such as yearly stock price movements, with one price per day. Assume that the distance function between two such sequences S and Q is the Euclidean distance

$$\mathcal{D}(S,Q) \equiv \left(\sum_{i=1}^{N} (S[i] - Q[i])^2\right)^{1/2}$$
(7.2)

where S[i] stands for the value of stock S on the *i*-th day. Clearly, computing the distance of two stocks will take 365 subtractions and 365 squarings in our example.

The idea behind the 'quick-and-dirty' test is to characterize a sequence with a single number, which will help us discard many non-qualifying sequences. Such a number could be, eg., the average stock price over the year: Clearly, if two stocks differ in their averages by a large margin, it is impossible that they will be similar. The converse is not true, which is exactly the reason we may have false alarms. Numbers that contain some information about a sequence (or a multimedia object, in general), will be referred to as 'features' for the rest of this paper. Using a good feature (like the 'average', in the stock-prices example), we can have a quick test, which will discard many stocks with a single numerical comparison for each sequence, a big gain over the 365 subtractions and squarings that the original distance function requires.

If using one feature is good, using two or more features might be even better, because they may reduce the number of false alarms (at the cost of making the 'quick-and-dirty' test a bit more elaborate and expensive). In our stock-prices example, additional features might be, eg., the standard deviation, or, even better, some of the discrete Fourier transform (DFT) coefficients, as we shall see in Chapter 8.

The end result of using k features for each of our objects is that we can map each object into a point in k-dimensional space. We shall refer to this mapping as $\mathcal{F}()$ (for 'F'eature):

Definition 7.2 Let $\mathcal{F}()$ be the mapping of objects to k-d points, that is $\mathcal{F}(O)$ will be the k-d point that corresponds to object O.

This mapping provides the key to improve on the second drawback of sequential scanning: by organizing these k-d points into a spatial access method, we can cluster them in a hierarchical structure, like the R-trees. Upon a query, we can exploit the R-tree, to prune out large portions of the database that are not promising. Thus, we do not even *have* to do the quick-and-dirty test on all of the k-d points!



Figure 7.1 Illustration of basic idea: a database of sequences S1, ... Sn; each sequence is mapped to a point in feature space; a query with tolerance ϵ becomes a sphere of radius ϵ .

Figure 7.1 illustrates the basic idea: Objects (eg., time sequences that are 365-points long) are mapped into 2-d points (eg., using the average and the standard-deviation as features). Consider the 'whole match' query that requires all the objects that are similar to S_n within tolerance ϵ : this query becomes an k-d sphere in feature space, centered on the image $\mathcal{F}(S_n)$ of S_n . Such queries on multidimensional points is exactly what R-trees and other SAMs are designed to answer efficiently. More specifically, the search algorithm for a whole match query is illustrated in Figure 7.2.

Algorithm 7.1 Search for whole-match queries:

- 1. map the query object Q into a point $\mathcal{F}(Q)$ in feature space
- 2. using the SAM, retrieve all points within the desired tolerance ϵ from $\mathcal{F}(Q)$.
- 3. retrieve the corresponding objects, compute their actual distance from Q and discard the false alarms.

Figure 7.2 Pseudo-code for the search algorithm.

Intuitively, this approach has the potential to relieve both problems of the sequential scan, presumably resulting into much faster searches. The only step that we have to be careful with is that the mapping $\mathcal{F}()$ from objects to k-d points does not distort the distances. Let $\mathcal{D}()$ be the distance function of two objects, and $\mathcal{D}_{feature}()$ be the (say, Euclidean) distance of the corresponding feature vectors. Ideally, the mapping should preserve the distances exactly, in which case the SAM will have neither false alarms nor false dismissals. However, requiring perfect distance preservation might be difficult: For example, it is not obvious which features we have to use to match the editing distance between two DNA strings. Even if the features are obvious, there might be practical problems: for example, in the stock-price example, we could treat every sequence as a 365-dimensional vector; although in theory a SAM can support an arbitrary number of dimensions, in practice all SAMs suffer from the 'dimensionality curse', as discussed in Chapter 5.

The crucial observation is that we can guarantee that the proposed method will not result in any false dismissals, if the distance in feature space matches or underestimates the distance between two objects. Intuitively, this means that our mapping $\mathcal{F}()$ from objects to points *should make things look closer*, i.e., it should be a contractive mapping.

Mathematically, let O_A and O_B be two objects (e.g., same-length sequences) with distance function $\mathcal{D}()$ (e.g., the Euclidean distance) and $\mathcal{F}(O_1)$, $\mathcal{F}(O_2)$ be their feature vectors (e.g., their first few Fourier coefficients), with distance function $\mathcal{D}_{feature}()$ (e.g., the Euclidean distance, again). Then we have:

Lemma 1 (Lower-bounding) To guarantee no false dismissals for wholematch queries, the feature extraction function $\mathcal{F}()$ should satisfy the following formula:

$$\mathcal{D}_{feature}(\mathcal{F}(O_1), \mathcal{F}(O_2)) \le \mathcal{D}(O_1, O_2) \tag{7.3}$$

for every pair of objects O_1 , O_2 .

Proof: Let Q be the query object, O be a qualifying object, and ϵ be the tolerance. We want to prove that if the object O qualifies for the query, then it will be retrieved when we issue a range query on the feature space. That is, we want to prove that

$$\mathcal{D}(Q, O) \le \epsilon \Rightarrow \mathcal{D}_{feature}(\mathcal{F}(Q), \mathcal{F}(O)) \le \epsilon$$

However, this is obvious, since

$$\mathcal{D}_{feature}(\mathcal{F}(Q), \mathcal{F}(O)) \le \mathcal{D}(Q, O) \le \epsilon$$

Thus, the proof is complete.

We have just proved that lower-bounding the distance works correctly for range queries. Will it work for the other queries of interest, like 'all-pairs' and 'nearest neighbor' ones? The answer is affirmative in both cases: An 'all-pairs' query can easily be handled by a 'spatial join' on the points of the feature space: using a similar reasoning as before, we see that the resulting set of pairs will be a superset of the qualifying pairs. For the nearest-neighbor query, the following algorithm guarantees no false dismissals: (a) find the point $\mathcal{F}(P)$ that is the nearest neighbor to the query point $\mathcal{F}(Q)$ (b) issue a range query, with query object Q and radius $\epsilon = \mathcal{D}(Q, P)$ (ie, the actual distance between the query object Q and data object P. For more details and for an application of this algorithm on tumor-like shapes, see [KSF⁺96].

In conclusion, the proposed generic approach to indexing multimedia objects for fast similarity searching shown in Figure 7.3 (named '*GEMINI*' for *GEneric Multimedia object INdexIng*):

The first two steps of GEMINI deserve some more discussion: The first step involves a domain expert. The methodology focuses on the *speed* of search only; the quality of the results is completely relying on the distance function that the expert will provide. Thus, GEMINI will return *exactly the same* response-set (and therefore, the same quality of output) with what the sequential scanning of the database would provide; the only difference is that GEMINI will be faster.

The second step of GEMINI requires intuition and imagination. It starts by trying to answer the question (referred to as the '*feature-extracting*' question for the rest of this work):

Algorithm 7.2 ('GEMINI') (GEneric Multimedia INdexIng approach):

- 1. determine the distance function $\mathcal{D}()$ between two objects
- 2. find one or more numerical feature-extraction functions, to provide a 'quick and dirty' test
- 3. prove that the distance in feature space lower-bounds the actual distance $\mathcal{D}()$, to guarantee correctness
- 4. use a SAM (eg., an R-tree), to store and retrieve the k-d feature vectors

Figure 7.3 Pseudo-code for the GEMINI algorithm.

'Feature-extracting' question: If we are allowed to use only one numerical feature to describe each data object, what should this feature be?

The successful answers to the above question should meet two goals: (a) they should facilitate step 3 (the distance lower-bounding) and (b) they should capture most of the characteristics of the objects.

We give case-studies of the GEMINI algorithm in the next two Chapters. The first involves 1-d time sequences, and the second focuses on 2-d color images. We shall see that the approach of the 'quick-and-dirty' filter, in conjunction with the lower-bounding lemma (Lemma 1), can lead to solutions to two problems:

- The dimensionality curse (time sequences).
- The 'cross-talk' of features (color images).

For each case study we (a) describe the objects and the distance function (b) show how to apply the lower-bounding lemma and (c) give experimental results, on real or realistic data. In Chapter 10 we show how to extend the idea of a 'quick-and-dirty' filter to handle sub-pattern matching in time sequences. In Chapter 11 we present '*FastMap*', an automated method of extracting features, for a given set of objects \mathcal{O} and for a given distance function $\mathcal{D}()$.

8

1-D TIME SEQUENCES

Here the goal is to search a collection of (equal-length) time sequences, to find the ones that are similar to a desirable sequence. For example, 'in a collection of yearly stock-price movements, find the ones that are similar to IBM'.

8.1 DISTANCE FUNCTION

According to GEMINI (Algorithm 7.2), the first step is to determine the distance measure between two time sequences. As in [AFS93], we chose the Euclidean distance (Eq. 7.2), because it is the distance of choice in financial and forecasting applications (e.g., [LeB92]). Fast indexing for additional, more elaborate distance functions that include time-warping [SK83] [WTK86] [RJ93] is the topic of ongoing research (eg., [GK95, JMM95]).

8.2 FEATURE EXTRACTION AND LOWER-BOUNDING

Having decided on the Euclidean distance as the dis-similarity measure, the next step of the GEMINI algorithm is to find some features that can lowerbound it. We would like a set of features that (a) preserve/lower-bound the distance and (b) carry much information about the corresponding time sequence, so that the false alarms are few. The second requirement suggests that we use 'good' features, that have much discriminatory power. In the stock-price example, a 'bad' feature would be, eg., the first-day's value: the reason is that two stocks might have similar first-day values, yet they may differ significantly from then on. Conversely, two otherwise similar sequences, may agree everywhere, except for the first day's values. At the other extreme, we could use the values of *all* 365 days as features. However, although this would perfectly match the actual distance, it would lead to the 'dimensionality curse' problem.

Clearly, we need some better features. Applying the second step of the GEMINI algorithm, we ask the 'feature-extracting' question: if we are allowed to use only one feature from each sequence, what would this feature be? A natural answer is the average. By the same token, additional features could be the average of the first half, of the second half, of the first quarter, etc. Or, in a more systematic way, we could use the coefficients of the Fourier transform, and, for our case, the Discrete Fourier Transform (DFT) (see, eg., [OS75], or Appendix B). For a signal $\vec{x} = [x_i], i = 0, \ldots, n-1$, let X_f denote the *n*-point DFT coefficient at the *f*-th frequency $(f = 0, \ldots, n-1)$. Also, let $\vec{X} = [X_f]$ be the *n*-point DFT transform of \vec{x} . Appendix B provides a quick introduction to the basic concepts of the DFT.

The third step of the GEMINI methodology is to show that the distance in feature space lower-bounds the actual distance. The solution is provided by Parseval's theorem [OS75], which states that the DFT preserves the energy of a signal, as well as the distances between two signals:

$$\mathcal{D}(\vec{x}, \vec{y}) = \mathcal{D}(\vec{X}, \vec{Y}) \tag{8.1}$$

where \vec{X} and \vec{Y} are Fourier transforms of \vec{x} and \vec{y} respectively.

Thus, if we keep the first $k \leq n$ coefficients of the DFT as the features, we lower-bound the actual distance:

$$\mathcal{D}_{feature}(\mathcal{F}(\vec{x}), \mathcal{F}(\vec{y})) = \sum_{f=0}^{k-1} |X_f - Y_f|^2$$

$$\leq \sum_{f=0}^{n-1} |X_f - Y_f|^2 = \sum_{i=0}^{n-1} |x_i - y_i|^2 \equiv \mathcal{D}(\vec{x}, \vec{y})$$

because we ignore positive terms from the above Equation. Thus, there will be *no false dismissals*, according to Lemma 1.

Notice that our GEMINI approach can be applied with *any* orthonormal transform, such as, the Discrete Cosine transform (DCT) (see [Wal91] or Appendix B.4), the wavelet transform (see [RBC⁺92], [PTVF92], or Appendix C) etc., because they all preserve the distance between the original and the transformed space. In fact, our response time will improve with the ability of the transform to concentrate the *energy*: the fewer the coefficients that contain most of the energy, the more accurate our estimate for the actual distance, the fewer the false alarms, and the faster our response time. The energy of a signal is the sum of squares of its elements (see Definition A.7 in the Appendix). Thus, the performance results presented next are just pessimistic bounds; better transforms will achieve even better response times.

In addition to being readily available, (eg., in mathematical symbolic manipulation packages, like 'Mathematica', 'S', 'maple' etc.), the DFT concentrates the energy in the first few coefficients, for a large class of signals, the colored noises. These signals have a skewed energy spectrum that drops as $O(f^{-b})$. The energy spectrum or power spectrum of a signal is the square of the amplitude $|X_f|$, as a function of the the frequency f (see Appendix B).

- For b=2, we have the so-called random walks or brown noise, which model successfully stock movements and exchange rates (e.g., [Man77]). Our mathematical argument for keeping the first few Fourier coefficients agrees with the intuitive argument of the Dow Jones theory for stock price movement (see, for example, [EM66]). This theory tries to detect primary and secondary trends in the stock market movement, and ignores minor trends. Primary trends are defined as changes that are larger than 20%, typically lasting more than a year; secondary trends show 1/3-2/3 relative change over primary trends, with a typical duration of a few months; minor trends last roughly a week. From the above definitions, we conclude that primary and secondary trends correspond to strong, low frequency signals while minor trends correspond to weak, high frequency signals. Thus, the primary and secondary trends are exactly the ones that our method will automatically choose for indexing.
- With even more skewed spectrum (b > 2), we have the black noises [Sch91].
 Such signals model successfully, for example, the water level of rivers and the rainfall patterns as they vary over time [Man77].
- with b=1, we have the *pink noise*. Birkhoff's theory [Sch91] claims that 'interesting' signals, such as musical scores and other works of art, consist of 'pink noise', whose energy spectrum follows $O(f^{-1})$. The argument of the theory is that white noise with $O(f^0)$ energy spectrum is completely unpredictable, while brown noise with $O(f^{-2})$ energy spectrum is too predictable and therefore 'boring'. The energy spectrum of pink noise lies in-between.

As an illustration of a 'colored noise', Figure 8.1 plots the closing prices of the IBM stock, 8/30/93 - 4/20/94, along with the amplitude spectrum in linear and logarithmic scales. Notice how skewed the spectrum is, as well as how closely it is approximated by the theoretically expected 1/f line.



Figure 8.1 (a) Closing prices of IBM stock, 8/30/93 - 4/20/94, (b) its DFT amplitude spectrum in linear scales and (c) in logarithmic scales, along with the theoretically expected 1/f line.

Additional financial data sets with similar behavior are available from ftp: //sfi. santafe. edu/ pub/ Time-Series/ competition (Dataset 'C', with the Swiss franc - U.S. dollar exchange rate), or from http: //www. ai.mit. edu/ stocks.html, where there are stock prices for several companies.

In addition to 1-d signals (stock-price movements and exchange rates), it is believed that several families of higher dimensionality signals belong to the family of 'colored noises', with skewed spectrum. For example, 2-d signals, like photographs, typically exhibit a few strong coefficients in the lower spatial frequencies. The JPEG image compression standard [Wal91] exactly exploits this phenomenon, effectively ignoring the high-frequency components of the Discrete Cosine Transform, which is closely related to the Fourier transform (see Appendix B.4).

8.3 EXPERIMENTS

The above indexing method was implemented and compared against sequential scanning, on a collection of synthetic random walks. See [AFS93] for more details. As expected, the sequential scanning was outperformed.



Figure 8.2 Breakup of the execution time, for range queries: response time versus number of DFT coefficients.

An interesting point is how to determine the number k of DFT coefficients to retain. Figure 8.2 shows the break-up of the response time versus the number k of DFT coefficients retained. The diamonds, squares and crosses indicate total time, post-processing time and R-tree time, respectively. Notice that, as we keep more coefficients, the R-tree becomes bigger and slower, but more accurate (fewer false alarms, and therefore shorter post-processing time). This trade-off reaches an equilibrium for k=2 or 3.

The major conclusions from the application of the GEMINI method on time sequences are the following:

- 1. GEMINI can be successfully applied to time sequences, and specifically to the ones that behave like 'colored noises' (stock prices movements, currency exchange rates, water-level in rivers etc.)
- 2. For signals with skewed spectrum like the above ones, the minimum in the response time is achieved for a small number of Fourier coefficients (k = 1 3). Moreover, the minimum is rather flat, which implies that a sub-optimal choice for k will give search time that is close to the minimum. Thus, with the help of the lower-bounding lemma and the energy-concentrating properties of the DFT, we managed to avoid the 'dimensionality curse'.

3. The success in 1-d sequences suggests that the proposed GEMINI method is promising for 2-d or higher-dimensionality signals, if those signals also have skewed spectrum. The success of JPEG (that uses DCT) indicates that real images indeed have a skewed spectrum.

Exercises

Exercise 8.1 [10] Write a program to generate a random walk. Let each step be the output of a fair coin tossing: +1 or -1, with probability 50% each.

Exercise 8.2 [10] Compute and plot the spectrum of the above random walk, using some existing DFT package; also plot the 1/f line.

Exercise 8.3 [15] Use some time sequences from, eg., [Ton90, BJR94], and compute their DFT spectrum. List your observations.

Exercise 8.4 [25] Experiment with the 'energy concentrating' properties of DFT versus the DCT, for the sequences of the previous exercise: For each sequence, keep the k strongest coefficients of the DFT and the DCT transform; plot the squared error (sum of squares of omitted coefficients), as a function of k, for each of the two transforms. List your observations.

Exercise 8.5 [20] Write a program that will generate 'pink noise', that is, a signal with 1/f amplitude spectrum. (Hint: use the random walk above from Exercise 8.1, compute its DFT spectrum, divide each Fourier coefficient appropriately, and invert).

9

2-D COLOR IMAGES

Retrieving images by content attracts high and increasing interest [JN92, OS95, FSN^+95]. Queries on content may focus on color, texture, shape, position, etc. Potential applications include medical image databases ('Give me other images that contain a tumor with a texture like this one'), photo-journalism ('Find images that have blue at the top and red at the bottom'), art, fashion, cataloging, retailing etc..

In this chapter we present a color indexing algorithm and specifically the distance functions and the application of the GEMINI approach. We focus on 'whole-match' queries, or, equivalently, 'queries by example', where the user chooses an existing image (or draws one with a sketch-pad) and asks for similar color images. This color indexing algorithm was developed within the QBIC system of IBM [FBF+94]; see [NBE+93, Equ93, FSN+95] for more details on the shape and texture indexing algorithms of QBIC and their precision-recall performance evaluation.

Past work on image retrieval has resulted in many systems and methods to match images according to color, shape, texture and relative position. For color, a typical choice is the color histograms which we describe next; for shape-matching, the turning angle [Hor86], the moments of inertia [FBF⁺94] or the pattern spectrum [KSF⁺96] are among the choices; for texture, the directionality, granularity and contrast [Equ93] are a good set of features; for the relative position, the 2-D strings method and its variants have been used [CSY87, PO95]. However, although there is a lot of work on image matching in the Machine Vision community and on fast searching in the database community, the former typically focuses on the quality of the matching, while the latter focuses on the speed of the search; it is only recently that the two communities have started collaborating [JN92, NBE+93], in an attempt to provide fast *and* effective image retrieval by content.

Recent surveys on image matching are in [FBF⁺94, PF96]. A presentation of image retrieval systems is in the special issue of the IEEE Computer (Sept. 1995) [GR95].

9.1 DISTANCE FUNCTION

We mainly focus on the color features, because color presents an interesting problem (namely, the 'cross-talk' of features), which can be resolved by the proposed 'GEMINI' approach (algorithm 7.2). Features for shape and texture are described in [FSN+95], and can be easily mapped into k-d points.

Each object is a color image, that is, a 2-d array of pixels; every pixel has 3 color components (eg., 1 byte for 'red', 1 byte for 'green' and 1 byte for 'blue'). For each image, we compute an *h*-element color histogram using *h* colors. Conceptually, *h* can be as high as 2^{24} colors, with each color being denoted by a point in a 3-dimensional color space. In practice, we cluster similar colors together using an agglomerative clustering technique [DH73a], and choose one representative color for each bucket (= 'color bin'). Typical numbers of color bins are h = 256 and h = 64. Each component in the color histogram is the percentage of pixels that are most similar to that color. Figure 9.1 gives an example of such a histogram of a fictitious photograph of a sunset: there are many red, pink, orange and purple pixels, but only few white and green ones.



Figure 9.1 An example of a color histogram of a fictitious sunset photograph: Many red, pink, orange, purple and blue-ish pixels; few yellow, white and greenish ones

Once these histograms are computed, one method to measure the distance between two histograms $(h \times 1 \text{ vectors}) \vec{x}$ and \vec{y} is given by

$$d_{hist}^{2}(\vec{x}, \vec{y}) = (\vec{x} - \vec{y})^{t} \times \mathbf{A} \times (\vec{x} - \vec{y}) = \sum_{i}^{h} \sum_{j}^{h} a_{ij} (x_{i} - y_{i}) (x_{j} - y_{j})$$
(9.1)

where the superscript 't' indicates matrix transposition, ' \times ' indicates matrix multiplication, and the color-to-color similarity matrix **A** has entries a_{ij} which describe the similarity between color *i* and color *j*, with $a_{ii} = 1$ for every *i*.

9.2 LOWER-BOUNDING

If we try to use the color-histograms as feature vectors in the GEMINI approach, there are two obstacles: (a) The 'dimensionality curse' (h may be large, e.g. 64 or 256 for color features) and, most importantly, (b) the quadratic nature of the distance function: The distance function in the feature space involves 'crosstalk' among the features (see Eq. 9.1), and it is thus a full quadratic form involving all cross terms. Not only is such a function much more expensive to compute than the Euclidean and any L_p distance, but it also precludes the use of spatial access methods ('SAMs'), because SAMs implicitly assume that there is no cross-talk. Figure 9.2 illustrates the situation: to compute the distance between the two color histograms \vec{x} and \vec{q} , the, eg., bright-red component of \vec{x} has to be compared not only to the bright-red component of \vec{q} , but also to the pink, orange etc. components of \vec{q} .



Figure 9.2 Illustration of the 'cross-talk' between two color histograms

To resolve the cross-talk problem, we resort to the 'GEMINI' approach (algorithm 7.2). The first step of the algorithm has been done: the distance function

between two color images is given by Eq. 9.1: $\mathcal{D}() = d_{hist}()$. The second step is to find one or more numerical features, whose Euclidean distance would lower-bound $d_{hist}()$. Thus, we ask the 'feature-extracting' question again: If we are allowed to use only one numerical feature to describe each color image, what should this feature be? According to the previous chapter on time sequences, we can consider some average value, or the first few coefficients of the 2-dimensional DFT transform. Since we have three color components, (eg., Red, Green and Blue), we could consider the average amount of red, green and blue in a given color image.

Notice that different color spaces can be used, with absolutely no change in the indexing algorithms. Thus, we continue the discussion with the RGB color space. This means that the color of an individual pixel is described by the triplet (R,G,B) (for 'R'ed, 'G'reen, 'B'lue). The average color of an image $\bar{x} = (R_{avg}, G_{avg}, B_{avg})^t$, is defined in the obvious way, with

$$R_{avg} = (1/P) \sum_{p=1}^{P} R(p)$$

$$G_{avg} = (1/P) \sum_{p=1}^{P} G(p)$$

$$B_{avg} = (1/P) \sum_{p=1}^{P} B(p)$$

where P is the number of pixels in the image, and R(p), G(p), and B(p) are the red, green and blue components (intensities, typically in the range 0-255) respectively of the p-th pixel. Given the average color vectors \bar{x} and \bar{y} of two images, we define $d_{avg}()$ as the Euclidean distance between the 3-dimensional average color vectors,

$$d_{avg}^2(\bar{x},\bar{y}) = (\bar{x}-\bar{y})^t \times (\bar{x}-\bar{y}) = \sum_{i=1}^3 (x_i - y_i)^2$$
(9.2)

The third step of the GEMINI algorithm is to prove that our feature distance $\mathcal{D}_{feature}() \equiv d_{avg}()$ lower-bounds the actual distance $\mathcal{D}() \equiv d_{hist}()$. Indeed, this is true, as an application of the so-called *QDB*- 'Quadratic Distance Bounding' Theorem [FBF+94].

The result is that, given a color query, our retrieval proceeds by first filtering the set of images based on their average (R, G, B) color, then doing a final,

more accurate matching using their full h-element histogram. The resulting speedup is discussed next.

9.3 EXPERIMENTS

Experiments are reported in $[FBF^+94]$, on a database of N=924 color image histograms, each of h=256 colors, of assorted natural images. The proposed method requires from a fraction of a second up to ≈ 4 seconds, while sequential scanning with the color-histogram distance (Eq. 9.1) requires ≈ 10 seconds. The performance gap is expected to increase for larger databases.

Thus, the conclusions are the following:

- The 'GEMINI' approach (ie., the idea to extract some features for a quickand-dirty test) motivated a fast method, using the average RGB distance $(d_{avg}())$; it also motivated a strong theorem (the so-called *QDB*- 'Quadratic Distance Bounding' Theorem [FBF+94]) which guarantees the correctness in our case.
- In addition to resolving the cross-talk problem, 'GEMINI' solved the 'dimensionality curse' problem at no extra cost, requiring only k=3 features, as opposed to h=64 or 256 that d_{hist}() required.

10 SUB-PATTERN MATCHING

10.1 INTRODUCTION

Up to now, we have examined the 'whole-match' case. The goal in this chapter is to extend the 'GEMINI' approach of the 'quick-and-dirty' test, so that we can handle sub-pattern matching queries. We focus on 1-d time series, to illustrate the problem and the solution more clearly. Then, the problem is defined as follows:

- We are given a collection of N sequences of real numbers S_1 , S_2 , S_N , each one of potentially different length.
- The user specifies query subsequence Q of length Len(Q) (which may vary) and the tolerance ϵ , that is, the maximum acceptable dis-similarity (= distance).
- We want to find quickly all the sequences S_i ($1 \le i \le N$), along with the correct offsets p, such that the subsequence $S_i[p:p+Len(Q) 1]$ matches the query sequence: $\mathcal{D}(Q, S_i[p:p+Len(Q) 1]) \le \epsilon$.

As in Chapter 8, we use the Euclidean distance as the dis-similarity measure $\mathcal{D}()$. The brute-force solution is to examine sequentially every possible subsequence of the data sequences for a match. We shall refer to this method by 'SequentialScan' method. Next, we describe a method that uses a small space overhead, to achieve up to 2 orders of magnitudes savings over the 'SequentialScan' method [FRM94].

10.2 SKETCH OF THE APPROACH - 'ST-index'

Without loss of generality, we assume that the minimum query length is w, where $w (\geq 1)$ depends on the application. For example, in stock price databases, analysts are interested in weekly or monthly patterns because shorter patterns are susceptible to noise [EM66]. Notice that we never lose the ability to answer shorter than w queries, because we can always resort to sequential scanning.



Figure 10.1 Illustration of the way that trails are created in feature space

Generalizing the reasoning of the method for 'whole matching', we use a sliding window of size w and place it at every possible position (offset), on every data sequence. For each such placement of the window, we extract the features of the subsequence inside the window. Thus, a data sequence of length Len(S) is mapped to a trail in feature space, consisting of (Len(S) - w + 1) points: one point for each possible offset of the sliding window. Figure 10.1 gives an example of a trail: Consider the sequence S_1 , and assume that we keep the first k=2 features (eg, the amplitude of the first and second coefficient of the w-point DFT). When the window of length w is placed at offset=0 on S_1 , we obtain the first point of the trail; as the window slides over S_1 , we obtain the rest of the trail.

The straightforward way to index these trails would be to keep track of the individual points of each trail, storing them in a spatial access method. We call this method *'I-naive'* method, where 'I' stands for 'Index' (as opposed to sequential scanning). However, storing the individual points of the trail in an R-tree is inefficient, both in terms of space as well as search speed. The reason is that, almost every point in a data sequence will correspond to a point in the k-dimensional feature space, leading to an index with a 1:k increase in storage requirements. Moreover, the search performance will also suffer because the R-tree will become tall and slow. As shown in [FRM94], the *'I-naive'* method



Figure 10.2 Example of (a) dividing trails into sub-trails and MBRs, and (b) grouping of MBRs in larger ones.

ended up being almost twice as slow as the 'SequentialScan'! Thus, we want to improve the 'I-naive' method, by making the representation of the trails more compact.

Here is where the idea of a 'quick-and-dirty' test leads to a solution: Instead of laboriously keeping track of each and every point of a trail in feature space, we propose to exploit the fact that successive points of the trail will probably be similar, because the contents of the sliding window in nearby offsets will be similar. We propose to divide the trail of a given data sequence into sub-trails and represent each of them with its minimum bounding (hyper)-rectangle (MBR). Thus, instead of storing thousands of points of a given trail, we shall store only a few MBRs. More importantly, at the same time we still guarantee 'no false dismissals': when a query arrives, we shall retrieve all the MBRs that intersect the query region; thus, we shall retrieve all the qualifying sub-trails, plus some false alarms (sub-trails that do not intersect the query region, while their MBR does).

Figure 10.2(a) gives an illustration of the proposed approach. Two trails are drawn; the first curve, labeled C1 (in the north-west side), has been divided into three sub-trails (and MBRs), whereas the second one, labeled C2 (in the south-east side), has been divided in five sub-trails. Notice that it is possible that MBRs belonging to the same trail may overlap, as C2 illustrates.

Thus, the idea is to map a data sequence into a *set of rectangles* in feature space. This yields significant improvements with respect to space, as well as

with respect to response time, as we shall see in section 10.3. Each MBR corresponds to a whole sub-trail, that is, points in feature space that correspond to successive positionings of the sliding window on the data sequences.

These MBRs can be subsequently stored in a spatial access method, such as an R-tree (see Chapter 5). Figure 10.2(b) shows how the eight leaf-level MBRs of Figure 10.2(a) will be grouped to form two MBRs at the next higher level, assuming a fanout of 4 (i.e. at most 4 items per non-leaf node). Note that the higher-level MBRs may contain leaf-level MBRs from different data sequences. For example, in Figure 10.2(b) notice that the left-side MBR1 contains a part of the C2 curve.

This completes the sketch of the proposed index structure. We shall refer to it by `ST-index', for 'Sub-Trail index'. There are two questions that we have to answer, to complete the description of the method:

- Insertions: When a new data sequence is inserted, what is a good way to divide its trail in feature space into sub-trails? The idea [FRM94] is to use an adaptive, heuristic algorithm, which will break the trail into sub-trails, so that the resulting MBRs of the sub-trails have small volume (and, therefore, result in few disk accesses on search).
- Queries: How to handle queries, and especially the ones that are longer than w? Figure 10.3 shows how to handle queries of length w: the query sub-sequence is translated into a point Q in feature space; all the MBRs that are within radius ϵ from Q are retrieved. Searching for longer queries is handled by breaking the query pattern into pieces of length w, searching the R-tree for each piece, and then merging the results. Figure 10.4 illustrates the algorithm: the query sequence is broken into p=2 pieces, each of length w; each piece gives rise to a range query in feature space.

More details and the correctness proofs for the above algorithms are in [FRM94].

10.3 EXPERIMENTS

The above method was implemented and tested on real stock-price time sequences totaling 329,000 points [FRM94]. The proposed method achieved up to 100 times better response time, compared to the sequential scanning.



Figure 10.3 Illustration of the search algorithm for minimum-length queries. The query becomes a sphere of radius ϵ in feature space.



Figure 10.4 Illustration of the search algorithm for a longer query. The query is divided in p=2 pieces of length w each, giving rise to p range queries in feature space

The conclusion is that the idea of using a 'quick-and-dirty' filter pays off again. Every sequence is represented *coarsely* by a set of MBRs in feature space; despite the loss of information, these MBRs provide the basis for quick filtering, which eventually achieves large savings over the sequential scanning.

Exercises

Exercise 10.1 [40] Implement a system which can search a collection of stockprices for user-specified patterns. Assume that the minimum length is one week; Use artificially generated random walks (exercise 8.1), or real data, from http: //www.ai.mit.edu/ stocks.html.

11 FASTMAP

11.1 INTRODUCTION

In the previous chapters we saw the 'GEMINI' approach, which suggests that we rely on domain experts to derive k feature-extraction functions, thus mapping each object into a point in k-dimensional space. Then the problem is reduced to storing, retrieving and displaying k-dimensional points, for which there is a plethora of algorithms available.

However, it is not always easy to derive the above feature-extraction functions. Consider the case, eg., of typed English words, where the distance function is the editing distance (minimum number of insertion, deletions and substitutions to transform one string to the other). It is not clear which the features should be in this case. Similarly, in matching digitized voice excerpts, we typically have to do some time-warping [SK83], which makes it difficult to design featureextraction functions.

Overcoming these difficulties is exactly the motivation behind this chapter. Automatically mapping the objects into points in some k-d space provides two major benefits:

- It can accelerate the search time for queries. The reason is that we can employ highly fine-tuned Spatial Access Methods (SAMs), like the R^{*}trees [BKSS90] and the z-ordering [Ore86]. These methods provide fast searching for range queries as well as spatial joins [BKSS94].
- 2. It can help with visualization, clustering and data-mining: Plotting objects as points in k=2 or 3 dimensions can reveal much of the structure of

the dataset, such as the existence of major clusters, the general shape of the distribution (linear versus curvilinear versus Gaussian) etc.. These observations can provide powerful insights in formulating hypotheses and discovering rules.

Next, we shall use the following terminology:

Definition 11.1 The k-dimensional point P_i that corresponds to the object O_i , will be called 'the *image*' of object O_i . That is, $P_i = (x_{i,1}, x_{i,2}, \ldots, x_{i,k})$

Definition 11.2 The k-dimensional space containing the 'images' will be called *target space*.

Given the above, the problem is defined as follows (see Figure 11.1 for an illustration):

$\mathbf{Problem}$

Given N objects and distance information about them (eg., an $N \times N$ distance matrix, or simply the distance function $\mathcal{D}(*,*)$ between two objects) **find** N points in a k-dimensional space, **such that** the distances are maintained as well as possible.

A special case is the situation where the N objects are n-d vectors. Then, goal is to do dimensionality reduction, mapping them into points in k-d space, preserving the (Euclidean) distances as well as possible. In this case, the optimal solution is given by the Karhunen-Loeve ('K-L') transform, which is described in detail in Appendix D.1.

For the general case, we expect that the distance function $\mathcal{D}()$ is non-negative, symmetric and obeys the triangular inequality. In the 'target' (k-d) space, we typically use the Euclidean distance, because it is invariant under rotations. Alternative distance metrics could be any of the L_p metrics, like the L_1 ('cityblock' or 'Manhattan' distance).

The ideal mapping should fulfill the following requirements:

FastMap



Figure 11.1 (a) Six objects and their pair-wise distance information; (b) 3-d points that try to approximate the given distance information.

- 1. It should be fast to compute: O(N) or $O(N \log N)$, but not $O(N^2)$ or higher, because the cost will be prohibitive for large databases.
- 2. It should preserve distances, leading to small discrepancies (low 'stress' see (Eq. 11.1)).
- 3. It should provide a very fast algorithm to map a new object (eg., a query object) to its image. The algorithm should be O(1) or $O(\log N)$. This requirement is vital for 'queries-by-example'.

The outline of this chapter is as follows. In section 11.2 we present a brief survey of Multi-Dimensional Scaling (MDS). In section 11.3 we describe 'FastMap', a linear, approximate solution [FL95]. In section 11.4 we show the results of FastMap on some real data (document vectors). Finally, section 11.5 lists the conclusions.

11.2 MULTI-DIMENSIONAL SCALING (MDS)

Multi-dimensional scaling (MDS) is used to discover the underlying (spatial) structure of a set of data items from the (dis)similarity information among them. There are several variations, but the basic method (eg., see [KW78]) is described next. The method expects (a) a set of N items, (b) their pair-wise (dis)similarities and (c) the desirable dimensionality k. Then, the algorithm

will map each object to a point in a k dimensional space, to minimize the *stress* function:

$$stress = \sqrt{\frac{\sum_{i,j} (\hat{d}_{ij} - d_{ij})^2}{\sum_{i,j} d_{ij}^2}}$$
 (11.1)

where d_{ij} is the dissimilarity measure between object O_i and object O_j and \hat{d}_{ij} is the (Euclidean) distance between their 'images': points P_i and P_j . The 'stress' function gives the relative error that the distances in k-d space suffer from, on the average.

To achieve its goal, MDS starts with a guess and iteratively improves it, until no further improvement is possible. In its simplest version, the algorithm works roughly as follows: It originally assigns each item to a k-d point (eg., using some heuristic, or even at random). Then, it examines every point, computes the distances from the other N - 1 points and moves the point to minimize the discrepancy between the actual dissimilarities and the estimated k-d distances. Technically, MDS employs the method of 'steepest descent' to update the positions of the k-d points. Intuitively, it treats each pair-wise distance as a 'spring' between the two points; then, the algorithm tries to re-arrange the positions of the k-d points to minimize the 'stress' of the springs.

MDS has been used in numerous, diverse applications, to help visualize a set of objects, given their pair-wise similarities. However, for our applications, MDS suffers from two drawbacks:

- It requires $O(N^2)$ time, where N is the number of items. Thus, it is impractical for large datasets. In the applications that MDS has been used, the number of items was small (typically, N=10-100).
- Its use for fast retrieval is questionable: In the 'query-by-example' setting, the query item has to be mapped to a point in k-d space. MDS is not prepared for this operation: Given that the MDS algorithm is $O(N^2)$, an incremental algorithm to search/add a new item in the database would be O(N) at best. Thus, the complexity of answering a query would be as bad as sequential scanning.

The above two drawbacks are the motivation behind the next section.

FastMap

Symbols	Definitions.
N	Number of objects in database
k	dimensionality of 'target space'
$\mathcal{D}(*,*)$	the distance function between two objects
$\parallel \vec{x} \parallel$	the length (= L_2 norm) of vector \vec{x}
(AB)	the length of the line segment AB

Table 11.1 Summary of Symbols and Definitions

11.3 A FAST, APPROXIMATE ALTERNATIVE: FASTMAP

The goal is to find N points in k-d space, whose Euclidean distances will match the distances of a given $N \times N$ distance matrix. Table 11.1 lists the symbols and their definitions. The key idea is to pretend that objects are indeed points in some unknown, n-dimensional space, and to try to project these points on k mutually orthogonal directions. The challenge is to compute these projections from the distance matrix only, since it is the only input we have. For the rest of this discussion, an object will be treated as if it were a point in an n-d space, with unknown n.

The heart of the 'FastMap' method is to project the objects on a carefully selected 'line'. To do that, we choose two objects O_A and O_B (referred to as 'pivot objects' from now on), and consider the 'line' that passes through them in *n*-d space. The algorithm to choose pivot objects uses a linear, approximate heuristic, and is discussed later (see Figure 11.4).

The projections of the objects on that line are computed using the cosine law:

$$d_{b,i}^{2} = d_{a,i}^{2} + d_{a,b}^{2} - 2x_{i}d_{a,b}$$
(11.2)

See Figure 11.2 for an illustration. Eq. 11.2 can be solved for x_i , the first coordinate of object O_i :

$$x_{i} = \frac{d_{a,i}^{2} + d_{a,b}^{2} - d_{b,i}^{2}}{2d_{a,b}}$$
(11.3)

In the above equations, d_{ij} is a shorthand for the distance $\mathcal{D}(O_i, O_j)$ (for $i, j = 1, \ldots, N$). Notice that the computation of x_i only needs the distances between objects, which are given.



Figure 11.2 Illustration of the 'cosine law' - projection on the line $O_A O_B$.

Observe that, thanks to Eq. 11.3, we can map objects into points on a line, preserving some of the distance information: For example, if O_i is close to the pivot O_A , then x_i will be small. Thus, we have the solution to our original problem, for k=1.



Figure 11.3 Projection on a hyper-plane \mathcal{H} , perpendicular to the line $O_A O_B$ of the previous figure.

To extend this method for 2-d and k-d target spaces, we keep on pretending that the objects are indeed points in n-d space: We consider a (n-1)-d hyperplane \mathcal{H} that is perpendicular to the line (O_A, O_B) ; then, project our objects on this hyper-plane. Let O_i' stand for the projection of O_i (for $i = 1, \ldots, N$).

FastMap

Algorithm 11.1 FastMap $(k, \mathcal{D}(), \mathcal{O})$

k: number of dimensions;

 $\mathcal{D}()$: the distance function;

 \mathcal{O} : the set of objects.

- 1. Pick-pivots O_A and O_B from \mathcal{O} .
- 2. Project all objects O_i along the line $O_A O_B$.
- 3. Call FastMap $(k 1, \mathcal{D}'(), \mathcal{O})$, unless k = 0.

Figure 11.4 Pseudo-code for the FastMap algorithm.

The problem is the same as the original problem, with n and k decreased by one.

The only missing part is to determine the distance function $\mathcal{D}'()$ between two of the projections on the hyper-plane \mathcal{H} , such as, O_i' and O_j' . Once this is done, we can recursively apply the previous steps. Figure 11.3 depicts two objects O_i, O_j , and their projections O_i', O_j' on the \mathcal{H} hyper-plane.

The distance function $\mathcal{D}'()$ is given by:

$$(\mathcal{D}'(O_i', O_j'))^2 = (\mathcal{D}(O_i, O_j))^2 - (x_i - x_j)^2 \quad i, j = 1, \dots, N$$
(11.4)

using the Pythagorean theorem on the triangle $O_i CO_j$ (see Figure 11.3).

Ability to compute the distance $\mathcal{D}'()$ allows us to project on a second line, lying on the hyper-plane \mathcal{H} , and, therefore, orthogonal to the first line (O_A, O_B) by construction. Thus, we can solve the problem for a 2-d 'target' space. More importantly, we can apply the same steps recursively, k times, thus solving the problem for any k.

Figure 11.4 gives the pseudo-code for FastMap.

To pick the pivot objects, we use a linear-time heuristic, to choose a pair of objects that are far-away from each other. Notice that finding the exact solution requires a quadratic number of steps $(O(N^2))$, because we have to examine every possible pair at least once.

Algorithm 11.2 Pick-pivots ($\mathcal{O}, \mathcal{D}()$)

- 1. Chose arbitrarily an object; let it be the second pivot object O_B
- 2. Set $O_A = ($ the object that is farthest apart from $O_B)$
- 3. Set $O_B = ($ the object that is farthest apart from $O_A)$
- 4. Repeat the last two steps a fixed number of times
- 5. Report the objects O_A and O_B as the pivot objects

Figure 11.5 Heuristic to choose two distant objects for pivots.

11.4 CASE STUDY: DOCUMENT VECTORS AND INFORMATION RETRIEVAL.

Here we use the algorithm on an information retrieval application [SM83]. As we mentioned in section 6.5, in the vector space model, documents are represented as vectors in V-dimensional space, where V is the size of the vocabulary of the collection. For the English language, we can expect V to range from 2,000 up to and exceeding 100,000 (the vocabulary of every-day English, and the size of a very detailed dictionary, respectively [Pet80]). The coordinates of such vectors are called *term weights* and can be binary ('1' if the term appears in the document; '0' if not) or real-valued, with values increasing with the importance (eg., occurrence frequency) of the term in the document.

Consider two documents d_1 and d_2 , with vectors $\vec{u_1}$, $\vec{u_2}$ respectively. As we mentioned, the similarity between two documents is typically measured by the *cosine similarity* of their vectors [SM83]:

$$similarity(d_1, d_2) = \frac{\vec{u_1} \circ \vec{u_2}}{\| \vec{u_1} \| \| \vec{u_2} \|}$$

where 'o' stands for the inner product of two vectors and ||*|| stands for the length (=Euclidean norm) of the vector. Clearly the cosine similarity takes values between -1 and 1. Figure 11.6 gives an example. There, $\cos(\theta)$ is considered as the similarity of the two vectors $\vec{u_1}$ and $\vec{u_2}$. Intuitively, the cosine similarity projects all the document vectors on the unit hyper-sphere (see vectors $\vec{u_{1,0}}$ and $\vec{u_{2,0}}$ in the Figure) and measures the cosine of the angle of the projections.

In order to apply FastMap, we first need to define a distance function that decreases with increasing similarity. From Figure 11.6 it would make sense
FastMap



Figure 11.6 Two vectors $\vec{u_1}$, $\vec{u_2}$, their angle θ and the cosine similarity function $\cos(\theta)$

to use the length of the line segment AB: $(AB) = || \vec{u_{1,0}} - \vec{u_{2,0}} ||$. After trigonometric manipulations, the result is

$$\mathcal{D}(d_1, d_2) = 2 * \sin(\theta/2) = \sqrt{2 * (1 - \cos(\theta))} = \sqrt{2 * (1 - similarity(d_1, d_2))}$$
(11.5)

Notice that Eq. 11.5 defines a distance function (non-negative, symmetric, satisfying the triangular inequality) and that it decreases with increasing similarity.

Also notice that it allows us to respond to range queries: Suppose that the user wants all the documents d that are similar to the query document q:

$$similarity(d,q) \ge \delta$$

Thanks to Eq. 11.5, the requirement becomes

$$\mathcal{D}(d,q) \le \sqrt{2*(1-\delta)} \tag{11.6}$$

which eventually becomes a range query in our 'target' space and can be handled efficiently by any SAM.

In [FL95], one of the testbeds we used was a collection of 35 text documents in 7 groups, with 5 documents per group: Abstracts of computer science technical

reports (labeled as 'Abs'), reports about basketball games ('Bbr'), 'call for papers' for technical conferences ('Cal'), portions of the Bible (from the Gospel of Matthew) ('Mat'), cooking recipes ('Rec'), 'world news' (documents about the Middle East - October 1994) ('Wor'), and sale advertisements for computers and software ('Sal')

Figure 11.7 shows the results of FastMap on the above dataset, with k=3 dimensions. The Figure shows the 3-d scatter-plot, (a) in its entirety and (b) after zooming into the center, to highlight the clustering abilities of FastMap. Notice that the 7 classes are separated well, in only k=3 dimensions. Additional experiments, involving real and synthetic datasets, are in [FL95].



Figure 11.7 The DOCS dataset, after FastMap in k=3-d space (a) The whole collection (b) magnification of the dashed box.

11.5 CONCLUSIONS

FastMap is a linear algorithm that maps N objects into k-d points, with small distortion of the distances. Thus it provides an automated way to extract features, for a given dataset \mathcal{O} and a given distance function $\mathcal{D}()$.

Mapping objects into points (while preserving the distances well) is vital for fast searching using the 'GEMINI' approach. Moreover, such a mapping is useful for data-mining, cluster analysis and visualization of a multimedia dataset. Notice that FastMap is linear on the number of objects N, while Multidimensional Scaling (MDS) is quadratic, and thus impractical for large databases.

Exercises

Exercise 11.1 [30] Implement MDS; apply it on some distance matrices

Exercise 11.2 [30] Implement the string-editing distance function; compute the distance matrix for 100 English words (eg., from /usr/dict/words), and map them into 2-d points using the MDS package of Exercise 11.1

Exercise 11.3 [25] Implement FastMap and compare it against MDS, with respect to its speed and its stress, on a set of 100, 200, 500 English words from /usr/dict/words. List your observations.

12 CONCLUSIONS

We have presented a generic method (the '*GEMINI*' approach) to accelerate queries by content on image databases and, more general, on multimedia databases. Target queries are, eg., 'find images with a color distribution of a sunset photograph'; or, 'find companies whose stock-price moves similarly to a given company's stock'.

The method expects a distance function $\mathcal{D}()$ (given by domain experts), which should measure the dis-similarity between two images or objects O_A , O_B . We mainly focus on whole match queries (that is, queries by example, where the user specifies the ideal object and asks for all objects that are within distance ϵ from the ideal object). Extensions to other types of queries (nearest neighbors, 'all pairs' and sub-pattern match) are briefly discussed.

The 'GEMINI' approach combines two ideas:

- The first is to devise a 'quick and dirty' test, which will eliminate several non-qualifying objects. To achieve that, we should extract k numerical features from each object, which should somehow describe the object (for example, the first few DFT coefficients for a time sequence, or for a gray-scale image). The key question to ask is 'If we are allowed to use only one numerical feature to describe each data object, what should this feature be?'
- The second idea is to further accelerate the search, by organizing these k-dimensional points using state of the art spatial access methods ('SAMs') like the R-trees. These methods typically group neighboring points together, thus managing to discard large un-promising portions of the address space early.

The above two ideas achieve fast searching. We went further, and we considered the condition under which the above method will be not only fast, but also *correct*, in the sense that it will not miss any qualifying object (false alarms are acceptable, because they can be discarded, with the obvious way). The answer is the *lower-bounding* lemma, which intuitively states that the mapping $\mathcal{F}()$ of objects to k-d points should make things look closer.

The rest of the chapters describe how to apply the method for a variety of environments, like 1-d time sequences and 2-d color images. These environments are specifically chosen, because they give rise to the 'dimensionality-curse' and the 'cross-talk' problems, respectively. The approach of the 'quick-and-dirty' filter, together with the lower-bounding lemma, provided solutions to both cases. Experimental results on real or realistic data illustrated the speed-up that the 'GEMINI' approach provides.

In the last two Chapters focused on two related problems. In Chapter 10 we presented a method to handle sub-pattern matching in time sequences, by using again a 'quick-and-dirty' filter and a sliding window, to map each sequence into a trail in feature space; the trail is represented coarsely by a small number of minimum bounding (hyper-)rectangles, which are fed into a Spatial Access Method (SAM). Finally, in Chapter 11 we discussed 'FastMap', a linear, approximate algorithm which can derive features automatically, given a set of objects \mathcal{O} and a distance function $\mathcal{D}()$.

PART III

MATHEMATICAL TOOLBOX

A PRELIMINARIES

Before we start, we need some definitions from complex algebra and from linear algebra. Consider a complex number

$$c = a + jb = A\exp(j\phi)$$

where $j = \sqrt{-1}$ is the imaginary unit. Then, we have the following:

Definition A.1 The amplitude |c| is defined as $A \equiv |c| = \sqrt{a^2 + b^2}$

Definition A.2 The <u>phase</u> ϕ of the number c = a + jb is defined as $\phi = \arctan(b/a)$

Definition A.3 The conjugate c^* of c is defined as a - jb.

Definition A.4 The energy E(c) of c is defined as the square of the amplitude $(E(c) \equiv |c|^2 \equiv c \ c^*).$

From matrix algebra, we use the following notation and concepts. \vec{x} will be considered as a column vector. Eg.,

$$\vec{x} = \begin{bmatrix} 2\\1\\3 \end{bmatrix} \tag{A.1}$$

In general, lower-case letters with an arrow will denote column vectors. Capital, bold letters (eg., **A**) will denote matrices: Eg., $\mathbf{A} = [a_{i,j}]$ where i, j span the rows and columns, respectively.

Definition A.5 \mathbf{A}^{t} denotes the <u>transpose</u> of a matrix: $\mathbf{A}^{t} = [a_{j,i}]$ (notice the reversal of i and j). If the matrix \mathbf{A} has complex entries, then we want the so-called <u>hermitian</u> matrix $\mathbf{A}^{*} = [a_{j,i}^{*}]$.

Clearly, for real-valued matrices, the transpose matrix is the hermitian matrix.

Definition A.6 The norm $|| \vec{x} ||$ of a vector \vec{x} is its Euclidean norm (root of sum of squares).

Definition A.7 The energy $E(\vec{x})$ of a sequence (or vector) \vec{x} is defined as the sum of energies at every point of the sequence:

$$E(\vec{x}) \equiv || \vec{x} ||^2 \equiv \sum_{i=0}^{n-1} |x_i|^2$$
(A.2)

Obviously, the energy of a signal is the square of the Euclidean norm (\equiv length) $\parallel \vec{x} \parallel$ of the vector \vec{x} .

Definition A.8 The <u>inner</u> or <u>'dot' product</u> $\vec{x} \circ \vec{y}$ of two vectors \vec{x} and \vec{y} is defined as

$$\vec{x} \circ \vec{y} = \sum_{i} (x_i * y_i) = (\vec{x})^t \times \vec{y}$$
(A.3)

where `*' denotes scalar multiplication and $'\times '$ denotes matrix multiplication

Obviously, $\| \vec{x} \|^2 = \vec{x} \circ \vec{x}$. For complex-valued vectors, we use the hermitian instead of the transpose.

Definition A.9 Orthogonality: Two vectors \vec{x} and \vec{y} are <u>orthogonal</u> ($\vec{x} - \vec{y}$) iff $\vec{x} \circ \vec{y} = 0$

Definition A.10 A matrix $\mathbf{A} = [\vec{a_1}, \vec{a_2}, \ldots]$ is <u>column-orthonormal</u> iff its column vectors $\vec{a_i}$ are unit vectors and mutually orthogonal, that is

$$\vec{a_i} \circ \vec{a_j} = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{otherwise} \end{cases}$$
(A.4)

Preliminaries

The definition for row-orthonormal is symmetric. From the above definition, if ${\bf A}$ is a column-orthonormal matrix, then

$$\mathbf{A}^t \times \mathbf{A} = \mathbf{I} \tag{A.5}$$

where I is the identity matrix of the appropriate dimensions. Again, for complex-valued matrices, we use the hermitian instead of the transpose.

B FOURIER ANALYSIS

B.1 DEFINITIONS

The intuition behind the Fourier Transform (as well as the Discrete Time Fourier transform that we are examining) is based on Fourier's theorem, that every continuous function can be considered as a sum of sinusoidal functions. For the discrete case, which is the one of interest to us, the *n*-point *Discrete Fourier Transform* [OS75] of a signal $\vec{x} = [x_i], i = 0, ..., n - 1$ is defined to be a sequence \vec{X} of *n* complex numbers $X_f, f = 0, ..., n - 1$, given by

$$X_f = 1/\sqrt{n} \sum_{i=0}^{n-1} x_i \exp\left(-j2\pi f i/n\right) \quad f = 0, 1, \dots, n-1$$
 (B.1)

where j is the imaginary unit $(j \equiv \sqrt{-1})$.

Compactly,

$$\vec{r} \iff \vec{X}$$
 (B.2)

will denote a DFT pair. The signal \vec{x} can be recovered by the inverse transform:

$$x_i = 1/\sqrt{n} \sum_{f=0}^{n-1} X_f \exp\left(j2\pi f i/n\right) \quad i = 0, 1, \dots, n-1$$
(B.3)

 X_f is a complex number, with the exception of X_0 , which is a real if the signal \vec{x} is real. There are some minor discrepancies among books: some define $X_f = 1/n \sum_{i=0}^{n-1} \ldots$ or $X_f = \sum_{i=0}^{n-1} \ldots$ We have followed the definition in (Eq B.1), for it simplifies the upcoming Parseval's theorem (Eq B.5).

Recall that $\exp(j\phi) \equiv \cos(\phi) + j\sin(\phi)$. The intuition behind DFT is to decompose a signal into sine and cosine functions of several frequencies, multiples of the basic frequency 1/n. The reasons of its success is that certain operations, like filtering, noise removal, convolution and correlation, can be executed more conveniently in the frequency domain. Moreover, the frequency domain highlights some properties of the signals, such as periodicity.

It is instructive to envision the DFT as a matrix operation:

Observation B.1 Eq. B.1 can be re-written as

$$\vec{X} = \mathbf{A} \times \vec{x} \tag{B.4}$$

where $\mathbf{A} = [a_{i,f}]$ is an $n \times n$ matrix with

$$a_{i,f} = 1/\sqrt{n} \exp(-j2\pi f i/n)$$
 $i, f = 0, \dots, n-1$

Observation B.2 Notice that \mathbf{A} is column-orthonormal, that is, its column vectors are unit vectors, mutually orthogonal. It is also true that \mathbf{A} is row-orthonormal, since it is a square matrix [PTVF92, p. 60].

That is,

$$\mathbf{A}^* \times \mathbf{A} = \mathbf{A} \times \mathbf{A}^* = \mathbf{I}$$

where **I** is the $(n \times n)$ identity matrix and **A**^{*} is the conjugate-transpose (*'her-mitian'*) of **A**, that is $\mathbf{A}^* = [a^*_{t,i}]$

The above observation has a very useful, geometric interpretation: since the DFT corresponds to a matrix multiplication with \mathbf{A} of Eq. B.4, and since the matrix \mathbf{A} is orthonormal, the matrix \mathbf{A} effectively does a rotation (but no scaling) of the vector \vec{x} in n-dimensional complex space; as a rotation, it does not affect the length of the original vector, nor the Euclidean distance between any pair of points.

B.2 PROPERTIES OF DFT

Next we list the properties of DFT that are most useful for our applications.

Theorem 1 (Parseval) Let \vec{X} be the Discrete Fourier Transform of the sequence \vec{x} . Then we have

$$\sum_{i=0}^{n-1} |x_i|^2 = \sum_{f=0}^{n-1} |X_f|^2$$
(B.5)

Proof: See, eg., [OS75].

An easier proof is based on Observation B.2 above. Intuitively, Parsevals' theorem states that the DFT preserves the energy (= square of the length) of the signal.

Property B.1 The DFT also preserves the Euclidean distance.

Proof: Using the fact that the DFT is equivalent to matrix multiplication (Eq. B.4), and that the matrix **A** is column-orthonormal, we can prove that the DFT also preserves the distance between two signals \vec{x} and \vec{y} . Let \vec{X} and \vec{Y} denote their Fourier transforms. Then, we have:

$$\| \vec{X} - \vec{Y} \|^{2} = \| \mathbf{A} \times \vec{x} - \mathbf{A} \times \vec{y} \|^{2}$$

$$= (\mathbf{A} \times \vec{x} - \mathbf{A} \times \vec{y})^{*} \times (\mathbf{A} \times \vec{x} - \mathbf{A} \times \vec{y})$$

$$= (\vec{x} - \vec{y})^{*} \times \mathbf{A}^{*} \times \mathbf{A} \times (\vec{x} - \vec{y})$$

$$= (\vec{x} - \vec{y})^{*} \times \mathbf{I} \times (\vec{x} - \vec{y})$$

$$= \| \vec{x} - \vec{y} \|^{2}$$

$$(B.6)$$

That is, the DFT maintains the Euclidean distance between the two signals \vec{x} and \vec{y} .

Observation B.2 gives a strong intuitive 'proof' for the above two properties. The geometric point of view of Observation B.2 is important: any transformation that corresponds to an orthonormal matrix **A** will also enjoy a theorem similar to Parseval's theorem for the DFT. Such transformations are the DCT and the DWT, that we examine later.

Property B.2 A shift in the time domain changes only the phase of the DFT coefficients, but not the amplitude:

$$[x_{i-i_0}] \Longleftrightarrow [X_f \exp\left(-2\pi f i_0 j/n\right)] \tag{B.7}$$

where ' \iff ' denotes a DFT pair. This is a useful property, if we are looking for, eg., matching stock prices with time-lags.

Property B.3 For real signals, we have $X_f = X_{n-f}^*$, for f = 1, 2, ..., n-1.

Proof: See [PTVF92, p. 511].

Using the above property, we only need to plot the amplitudes up to the middle, and specifically, up to q, if n = 2q + 1, or q + 1, if the duration is n = 2q.

Definition B.1 The resulting plot of $|X_f|$ versus f will be called the <u>amplitude</u> <u>spectrum</u> or plain <u>spectrum</u> of the given time sequence; its square will be the energy spectrum or power spectrum.

The <u>phase spectrum</u> is defined similarly, but it is not used as much as the amplitude and energy spectra.

Property B.4 The DFT requires $O(n \log n)$ computation time.

Although a straightforward computation of the DFT coefficients requires $O(n^2)$ time, the celebrated *Fast Fourier Transform* (FFT) exploits regularities of the $e^{j2\pi f i/n}$ function, and achieves $O(n \log n)$ time (eg., see [PTVF92]).

B.3 EXAMPLES

The main point we want to illustrate is the ability of the DFT to highlight the periodicities of the input signal \vec{x} . This is achieved through the *amplitude spectrum* that we defined before. Next, we give some carefully selected signals and their spectra.

Example B.1 A composite tone:

$$x_i = 6\sin(2\pi 4i/n + 0.5) + 3\sin(2\pi 8i/n) \quad i = 0, \dots, 31$$
 (B.8)



Figure B.1 A composite tone and its amplitude spectrum. Notice the spikes for the two component frequencies at f=4 and 8.

This is a sum of sinusoidal functions of frequencies 4 and 8. Digitized voice and musical sounds are sums of a few sinusoidal functions. Figure B.1 shows the signal in the time domain, and its amplitude spectrum. In the latter, notice the two clear peaks, at frequencies 4 and 8, with amplitudes proportional to 6 and 3, respectively, with a constant of proportionality: $\sqrt{n}/2$.



Figure B.2 A (shifted) impulse function, and its amplitude spectrum. Notice the 'frequency leak' in all the frequencies.

Example B.2 The 'impulse function', or 'Dirac delta' function: $x_0 = 1$; $x_i = 0$ for i > 0.

Figure B.2 shows an impulse function (shifted to the right, for better plotting), along with its spectrum. Notice that there is no dominating frequency in the amplitude spectrum. This is expected, since the original signal has no periodicities. This phenomenon is informally called *frequency leak*: the given signal has strong components on every frequency, and thus it can not be approximated (compressed) well by using few DFT coefficients.

In general, spikes and discontinuities require all the frequencies. Since the DFT effectively assumes that the signal is repeated infinite times, periodically (with period n), a high difference between x_0 and x_{n-1} also leads to frequency leak. This is illustrated with the next example:



Figure B.3 The ramp function, and its amplitude spectrum. Notice the frequency leak, again.

Example B.3 The 'ramp' function:

 $x_i = i$

Figure B.3 shows the the ramp function and its DFT. Notice that it also has a 'frequency leak', having non-zero amplitudes for all its DFT coefficients. However, the amplitudes are decreasing, compared to the impulse function.

This is an important observation: in general, if the input signal has a trend, the DFT has a frequency leak. Notice that the upcoming DCT avoids this specific problem.

Example B.4 Number of trappings for Canadian lynx (animals per year, 1821-1934).

See Figure B.4. This is a well-known dataset in population ecology [Sig93, p. 45], as well as time-sequence analysis [Ton90, BJR94]. It has a strong periodical nature, which is highlighted by the spike in the spectrum. It is interesting to notice that the population of hares (not shown here) also follows a similar periodic pattern, with the same period, but with some phase lead, because hares are a major food source for the Canadian lynx.



Figure B.4 Canadian lynx trappings (1821-1934), and its amplitude spectrum. Notice the spike at f=12, corresponding to a period of 9.5 years.

Figure B.5 highlights the ability of the DFT to concentrate the energy. Consider the spectrum of Figure B.4(b) and set to zero all the amplitudes, except for the two strongest ones (for f=0 and 12); this is the spectrum illustrated in Figure B.5(b). Figure B.5(a) shows the corresponding sequence in the time domain (by doing the inverse DFT), as well as the original 'lynx' dataset. Notice how well the approximate sequence matches the original, despite the tiny number of coefficients kept.

B.4 DISCRETE COSINE TRANSFORM (DCT)

For our purposes, the ideal transform should concentrate the energy into as few coefficients as possible, for most of the signals of interest. For several real signals, successive values are correlated: eg., in images, if a pixel is dark, chances are that its neighbors will also be dark. In these cases, the Discrete Cosine Transform (DCT) achieves better energy concentration than the DFT, and very close to optimal [Gal91, p. 54].



Figure B.5 (a) The original lynx dataset (with 'diamonds') and its approximation (with 'crosses') (b) the spectrum of the approximation.

Moreover, the DCT avoids 'frequency leak' problems that plague the DFT when the input signal has a 'trend' (see Example B.3). The DCT solves this problem cleverly, by conceptually reflecting the original sequence in the time axis around the last point and taking the DFT on the resulting, twice-as-long sequence. Exactly because the (twice-as-long) signal is symmetric, all the coefficients will be *real* numbers. Moreover, from the property B.3 of the DFT, their amplitudes will be symmetric along the middle $(X_f = X_{2n-f})$. Thus, we need to keep only the first n of them.

The formulas for the DCT are

$$X_f = 1/\sqrt{n} \sum_{i=0}^{n-1} x_i \cos \frac{\pi f(i+0.5)}{n} \quad f = 0, \dots, n-1$$
 (B.9)

and for the inverse:

$$x_i = 1/\sqrt{n}X_0 + 2/\sqrt{n}\sum_{f=1}^{n-1} X_f \cos \frac{\pi f(i+0.5)}{n} \quad i = 0, \dots, n-1$$
 (B.10)

As with the DFT, the complexity of the DCT is also $O(n \log(n))$.

B.5 *m*-DIMENSIONAL DFT/DCT (JPEG)

All the above transforms can be extended to *m*-dimensional signals: for m=2 we have gray-scale images, for m=3 we have 3-d MRI brain scans etc. Informally,

we have to do the transformation along each dimension: for example, for the DFT for a 1024x1024 matrix (eg., image), we have to do the DFT on each row, and then do the DFT on each column. Formally, for an $n_1 \times n_2$ array $[x_{i_1,i_2}]$ these operations are expressed as follows:

$$X_{f_1,f_2} = \frac{1}{\sqrt{n_1}} \frac{1}{\sqrt{n_2}} \sum_{i_1=0}^{n_1} \sum_{i_2=0}^{n_2} x_{i_1,i_2} \exp\left(-2\pi j i_1 f_1/n_1\right) \exp\left(-2\pi j i_2 f_2/n_2\right)$$

where x_{i_1,i_2} is the value (eg., gray scale) of the position (i_1, i_2) of the array, and f_1 , f_2 are the spatial frequencies, ranging from 0 to (n_1-1) and (n_2-1) respectively.

The formulas for higher dimensionalities m are straightforward. The formulas for the m-d inverse DFT and DCT are analogous. Notice that the 2-dimensional DCT is used in the JPEG standard [Wal91, Jur92] for image and video compression.

B.6 CONCLUSIONS

We have discussed some powerful, classic tools from signal processing, namely the DFT and DCT. The DFT is helpful in highlighting periodicities in the input signal, through its amplitude spectrum. The DCT is closely related to the DFT, and it has some additional desirable properties: its coefficients are always real (as opposed to complex), it handles well signals with trends, and it is very close to the optimal for signals whose successive values are highly correlated.

C WAVELETS

C.1 MOTIVATION

The wavelet transform is believed to avoid the 'frequency leak' problem even better. Consider the case of an impulse function (Example B.2): both in the DFT and the DCT transform, it has non-zero amplitudes in all frequencies. Thus, what would take a single number to describe in the time domain, will require several numbers in the frequency domain. The problem is that the DFT has no temporal locality: each of its coefficients provide information about all the time instants. A partial remedy would be the so-called 'Short Window Fourier Transform' (SWFT) [RV91]: We can divide the time sequence into frames of, say, w consecutive (non-overlapping) samples, and do the w-point DFT in each of these windows. Thus, an impulse function in the time domain will have a restricted 'frequency leak'. Figure C.1 shows intuitively what happens: In the time domain, each value gives the full information about that instant (but no information about frequencies). The DFT has coefficients that give full information about a given frequency, but it needs all the frequencies to recover the value at a given instant in time. The SWFT is somewhere in between.

The only non-elegant point of the SWFT is the choice of the width w of the window: How large should it be, and why? The solution to this problem is very clever: let the width w be variable! This is the basis for the Discrete Wavelet Transform (DWT). Figure C.1 illustrates how the DWT coefficients tile the frequence-time plane.



Figure C.1 Tilings of the time-frequency plane: (a) Time-domain representation (b) Discrete Fourier transform (DFT) (c) Short-Window Fourier transform (SWFT) (d) Discrete Wavelet transform (DWT).

C.2 DESCRIPTION

Several Discrete Wavelet transforms that have been proposed. The simplest to describe and code is the *Haar* transform. *Ignoring* temporarily some proportionality constants, the Haar transform operates on the whole signal, giving the sum and the difference of the left and right part; then it focuses recursively on each of the halves, and computes the difference of their two sub-halves, etc, until it reaches an interval with one only sample in it.

It is instructive to consider the equivalent, bottom-up procedure. The input signal \vec{x} must have a length n that is a power of 2, by appropriate zero-padding if necessary.

- 1. Level 0: take the first two sample points x_0 and x_1 , and compute their sum $s_{0,0}$ and difference $d_{0,0}$; do the same for all the other pairs of points (x_{2i}, x_{2i+1}) . Thus, $s_{0,i} = C * (x_{2i} + x_{2i+1})$ and $d_{0,i} = C * (x_{2i} - x_{2i+1})$, where C is a proportionality constant, to be discussed soon. The values $s_{0,i}$ ($0 \le i \le n/2$) constitute a 'smooth' (=low frequency) version of the signal, while the values $d_{0,i}$ represent the high-frequency content of it.
- 2. Level 1: consider the 'smooth' $s_{0,i}$ values; repeat the previous step for them, giving the even-smoother version of the signal $s_{1,i}$ and the smooth-differences $d_{1,i}$ $(0 \le i \le n/4)$
- 3. ... and so on recursively, until we have a smooth signal of length 2.

The Haar transform of the original signal \vec{x} is the collection of all the 'difference' values $d_{l,i}$ at every level l and offset i, plus the smooth component $s_{L,0}$ at the last level L ($L = \log_2(n) - 1$).

Wavelets

Following the literature, the appropriate value for the constant C is $1/\sqrt{2}$, because it makes the transformation matrix to be orthonormal (eg., see Eq. C.4). Adapting the notation (eg., from [Cra94] [VM]), the Haar transform is defined as follows:

$$d_{l,i} = 1/\sqrt{2} \, \left(s_{l-1,2i} - s_{l-1,2i+1} \right) \quad l = 0, \dots, L, \quad i = 0, \dots, n/2^{l+1} - 1 \quad (C.1)$$

with

$$s_{l,i} = 1/\sqrt{2} (s_{l-1,2i} + s_{l-1,2i+1})$$
 $l = 0, \dots, L, i = 0, \dots, n/2^{l+1} - 1$ (C.2)

with the initial condition:

$$s_{-1,i} = x_i \tag{C.3}$$

For example, the 4-point Haar transform is as follows. Envisioning the input signal \vec{x} as a column vector, and its Haar transform \vec{w} as another column vector ($\vec{w} = [s_{1,0}, d_{1,0}, d_{0,0}, d_{0,1}]^t$), the Haar transform is equivalent to a matrix multiplication, as follows:

$$\begin{bmatrix} s_{1,0} \\ d_{1,0} \\ d_{0,0} \\ d_{0,1} \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 & 1/2 & 1/2 \\ 1/2 & 1/2 & -1/2 & -1/2 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 & 0 \\ 0 & 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \times \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$
(C.4)

The above procedure is shared among *all* the wavelet transforms: We start at the lowest level, applying two functions at successive windows of the signal: the first function does some smoothing, like a weighted average, while the second function does a weighted differencing; the smooth (and shortened) version of the signal is recursively fed back into the loop, until the resulting signal is too short.

The Haar transform is still criticized for 'frequency leak' problems [Dau92, p. 10]. One of the most popular wavelet transforms is the so-called Daubechies-4 [Dau92]. We describe the derivation of the 0-th level of coefficients only, because the rest of the levels are derived recursively, as explained above. At each step, the Daubechies-4 DWT operates on 4 consecutive sample points; the 'smooth' component is given by

$$s_{0,i} = h_0 x_{2i} + h_1 x_{2i+1} + h_2 x_{2i+2} + h_3 x_{2i+3} \quad i = 0, \dots, n/2$$
(C.5)

and the 'difference' component is given by

$$d_{0,i} = h_3 x_{2i} - h_2 x_{2i+1} + h_1 x_{2i+2} - h_0 x_{2i+3} \quad i = 0, \dots, n/2$$
(C.6)

where

$$h_0 = \frac{1+\sqrt{3}}{4\sqrt{2}} \quad h_1 = \frac{3+\sqrt{3}}{4\sqrt{2}} \quad h_2 = \frac{3-\sqrt{3}}{4\sqrt{2}} \quad h_3 = \frac{1-\sqrt{3}}{4\sqrt{2}} \quad (C.7)$$

Notice that the signal is supposed to 'wrap-around' (ie., $x_{n+i} = x_i$ whenever the index *i* exceeds *n*). More details are in [PTVF92]. The code in **nawk** [AKW88] and Bourne 'shell' is attached in the next section.

Figure C.2 shows the basis functions of the Daubechies-4 DWT for n=32 points. The top left gives the basis function #5: it is the level-2 wavelet, starting at position 0. Notice that it mainly concentrates on the first half of the signal, giving a weighted difference. The top right gives the basis function #9: it is the level-1 wavelet starting at position 0. Notice that it has a shorter time-span than the previous (#5), but more violent oscillation (thus, higher frequency content). The bottom row shows the basis functions #17 and #18. They correspond to level-0 wavelets starting at offsets t=0 and t=2, respectively. As expected, the basis functions of this level have the shortest time-span and the highest frequency content. Also as expected, these two basis functions have identical shape and only differ by a horizontal shift.

C.3 DISCUSSION

The computational complexity of the above transforms is O(n), as it can be verified from Eq. C.1-C.3. Notice that this is faster than the $O(n \log(n))$ of FFT, without even the need to resort to any of the FFT-like techniques.

In addition to their computational speed, there is a fascinating relationship between wavelets, multiresolution methods (like quadtrees or the pyramid structures in machine vision), and fractals. The reason is that wavelets, like quadtrees, will need only a few non-zero coefficients for regions of the image (or the time sequence) that are smooth/homogeneous, while they will spend more effort on the 'high activity' areas. It is believed [Fie93] that the mammalian retina consists of neurons which are tuned each to a different wavelet. Naturally occurring scenes tend to excite only few of the neurons, implying that a wavelet transform will achieve excellent compression for such images. Similarly, the human ear seems to use a wavelet transform to analyze a sound, at least in the very first stage [Dau92, p. 6] [WS93].



Figure C.2 Basis functions #5, 9, 17, 18, for the Daubechies-4 DWT.

C.4 CODE FOR DAUBECHIES-4 DWT

Next we give the source code for the Daubechies-4 DWT. We have chosen **nawk** [AKW88] because (a) it is a language on a higher level than 'C' and (b) it is more widely available than the even higher level languages, of mathematical packages (like 'Mathematica', 'Maple', 'MatLab' etc). See [PTVF92, VM] for wavelet code in some of the above languages. Object code for several wavelet transforms is available in the 'xwpl' package (http://www.math.yale.edu:80 /wavelets/) from Yale.

#! /bin/sh -f

nawk '

```
# implements the Daubechies-4 wavelet -
# Following "Numerical Recipes in C", p. 593.
```

```
# Author: Christos Faloutsos, 1996
# Notice: the input signal MUST HAVE a length that is
#
    a power of 2 (and greater or equal to 4)
# Expected input format: a sequence of numbers,
# separated by white space (tabs, blanks, newlines)
BEGIN{ c[0] = (1+sqrt(3))/(4*sqrt(2));
      c[1] = (3+sqrt(3))/(4*sqrt(2));
      c[2] = (3-sqrt(3))/(4*sqrt(2));
      c[3] = (1-sqrt(3))/(4*sqrt(2));
      TOL = 10^{(-6)};
      count = 0;
       function abs(xx) {
          res = xx
          if (xx < 0) { res = -xx }
          return res
       }
       # chopArray
       function chopArray( xarg, Narg){
          for(ii=1; ii<=Narg; ii++){</pre>
              if( abs(xarg[ii]) < TOL) {xarg[ii] =0}</pre>
          }
       }
       # print array
       function printArray ( x, N){
          for(i=1; i<=N; i++){ printf "%g\n", x[i] }</pre>
       }
       *****
       # wraps the ivalue in the 1-Nval interval
       function wrap ( ival, Nval) {
          resval = ival-1;
          resval = resval % Nval
          resval ++
          return resval
       }
```

```
118
```

}

Wavelets

```
*****
# performs one step of the DWT transform
# on array xarg[1:Narg]
# Narg: should be a power of 2, and > 4
# It does the changes IN PLACE
******
function oneStepDWT ( xarg, Narg ) {
   jj = 0;
   for( ii=1; ii<Narg; ii +=2 ){</pre>
      jj ++;
      sres[jj] = c[0]*xarg[wrap(ii,Narg)] + \
              c[1]*xarg[wrap(ii+1,Narg)] + \
              c[2]*xarg[wrap(ii+2,Narg)]+ \
              c[3]*xarg[wrap(ii+3,Narg)];
      dres[jj] = c[3]*xarg[wrap(ii,Narg)] - \
              c[2]*xarg[wrap(ii+1,Narg)] + \
               c[1]*xarg[wrap(ii+2,Narg)] - \
              c[0]*xarg[wrap(ii+3,Narg)];
   }
   for( ii=1; ii<= Narg/2; ii++ ){</pre>
      xarg[ii] = sres[ii];
      xarg[ii + Narg/2 ] = dres[ii]
   }
   return
}
*****
# Does the full wavelet transform -
# it calls repeatedly the oneStepDWT()
# The array xarg[1,N] is changed IN PLACE
*****
function DWT(xarg, Narg){
   # assert that Narg >= 4 and Narg: power of 2
   # WILL NOT WORK OTHERWISE
   for( len=Narg; len>=4; len = len/2){
      oneStepDWT(xarg, len)
   }
}
******
```

read in the elements of the array

119

```
{ for( j = 1; j<= NF; j++) { count++; x[count] = $j } }
END {
    N =count; # array length
    DWT(x,N)
    chopArray(x,N)
    printArray(x,N)
}
' $*</pre>
```

C.5 CONCLUSIONS

The Discrete Wavelet Transform (DWT) achieves even better energy concentration than the DFT and DCT transforms, for natural signals [PTVF92, p. 604]. It uses multiresolution analysis, and it models well the early signal processing operations of the human eye and human ear.

K-L AND SVD

D.1 THE KARHUNEN-LOEVE (K-L) TRANSFORM

Before we examine the K-L transform, we should give the definition of eigenvalues and eigenvectors of a square matrix \mathbf{S} . (We use the letter 'S' to stress the fact that the matrix is square).

Definition D.1 For a square $n \times n$ matrix **S**, the unit vector \vec{x} and the scalar λ that satisfy

$$\mathbf{S} \times \vec{x} = \lambda \times \vec{x} \tag{D.1}$$

are called an eigenvector and its corresponding eigenvalue of the matrix S.

The eigenvectors of a *symmetric* matrix are mutually orthogonal and its eigenvalues are real. See [PTVF92, p. 457] for more details.

The intuitive meaning of these concepts is the following: A matrix **S** defines an affine transformation $\vec{y} = \mathbf{S} \times \vec{x}$, that involves rotation and/or scaling; the eigenvectors are the unit vectors along the directions that are *not* rotated by **S**; the corresponding eigenvalues show the scaling. For example, the matrix

$$\mathbf{S} = \begin{bmatrix} 2 & 1\\ 1 & 3 \end{bmatrix} \tag{D.2}$$

gives the ellipse of Figure D.1, when applied to the periphery of the unit circle. The major and minor axes of the ellipse correspond to the directions of the

(D.3)

two eigenvectors of S. Specifically, the strongest eigenvalue corresponds to the eigenvector of the major axis:



Figure D.1 The matrix S as a transformation: the unit circle becomes an ellipse, with major axis along the first eigenvector

The following observation will be used later, to show the strong connection between the eigenvalue analysis and the upcoming Singular Value Decomposition (Theorem 2).

Observation D.1 If **S** is a real and symmetric matrix (i.e., $\mathbf{S} = \mathbf{S}^t$), then it can be written in the form

$$\mathbf{S} = \mathbf{U} \times \mathbf{\Lambda} \times \mathbf{U}^t \tag{D.4}$$

where the columns of U are the eigenvectors of S and Λ is a diagonal matrix, with values the corresponding eigenvalues of S.

Notice that the above observation is a direct consequence of the definition of the eigenvalues and eigenvectors (Definition D.1) and the fact that the eigenvectors are mutually orthogonal, or, equivalently, that \mathbf{U} is column-orthonormal.

As an arithmetic example, for the matrix \mathbf{S} that we used in the example (see Eq. D.2), we have

$$\mathbf{S} = \begin{bmatrix} 2 & 1 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 0.52 & 0.85 \\ 0.85 & -0.52 \end{bmatrix} \times \begin{bmatrix} 3.62 & 0 \\ 0 & 1.38 \end{bmatrix} \times \begin{bmatrix} 0.52 & 0.85 \\ 0.85 & -0.52 \end{bmatrix}$$

D.1.1 K-L: Problem definition

Consider the following problem: Given a collection of n-d points, project them on a k-d sub-space $(k \ll n)$, minimizing the error of the projections (sum of squared differences). The problem has been studied extensively in statistical pattern recognition and matrix algebra. The optimal way to project ndimensional points onto k-dimensional points $(k \le n)$ is the Karhunen-Loève ('K-L') transform (eg., see [DH73b], [Fuk90]). In other words, the K-L transforms gives linear combination of axis (='attributes' = 'features'), sorted in 'goodness' order.

Figure D.2 gives an illustration of the problem and the solution: it shows a set of 2-d points, and the corresponding 2 directions (x' and y') that the K-L transform suggests: If we are allowed only k=1, the best direction to project on is the direction of x'; the next best is y' etc.



Figure D.2 Illustration of the Karhunen-Loève transformation - the 'best' axis to project is x'.

Next we give the detailed steps of the K-L, as well as the code in 'mathematica'. The K-L computes the eigenvectors of the *covariance* matrix (see Eq. D.7), sorts them in decreasing eigenvalue order, and approximates each data vector with its projections on the first k eigenvectors. The $n \times n$ covariance matrix **C** is defined as follows. Consider the $N \times n$ data matrix **A**, where rows correspond to data vectors and columns correspond to attributes. That is, $\mathbf{A} = [a_{ij}]$ (i = 1, ...N and j = 1, ...n). The covariance matrix $\mathbf{C} = [c_{pq}]$ roughly gives the attribute-to-attribute similarity, by computing the un-normalized correlation between the

two attributes p and q. Let $\overline{a_{..p}}$ be the average of the p-th column/attribute:

$$\overline{a_{..p}} = 1/N \sum_{i=1}^{N} a_{ip} \quad p = 1, 2, ...n$$
 (D.5)

Then, the entry $c_{p,q}$ of the covariance matrix ${\bf C}$ is

$$c_{p,q} = \sum_{i=1}^{N} (a_{i,p} - \overline{a_{i,p}}) (a_{i,q} - \overline{a_{i,q}})$$
(D.6)

Intuitively, we move the origin to the center of gravity of the given vectors, obtaining the matrix $\mathbf{B} = [b_{ij}] = [a_{ij} - \overline{a_{..j}}]$. We shall refer to the matrix \mathbf{B} as the 'zero-mean' matrix, exactly because its column averages are zero, by construction. Then, we compute the covariance matrix \mathbf{C} as follows:

$$\mathbf{C} = \mathbf{B}^t \times \mathbf{B} \tag{D.7}$$

Example D.1 Consider the data vectors $\vec{a_1} = [1, 2]^t$, $\vec{a_2} = [1, 1]^t$ and $\vec{a_3} = [0, 0]^t$. Then we have for the data matrix **A**:

$$\mathbf{A} = \begin{bmatrix} 1 & 2\\ 1 & 1\\ 0 & 0 \end{bmatrix}$$

The column averages are the coordinates of the center of gravity:

$$\overline{a_{\cdot,1}} = 2/3 \quad \overline{a_{\cdot,2}} = 1$$

The 'zero-mean' matrix \mathbf{B} is

$$\mathbf{B} = \begin{bmatrix} 1/3 & 1\\ 1/3 & 0\\ -2/3 & -1 \end{bmatrix}$$

and the covariance matrix ${\bf C}$ is

$$\mathbf{C} = \left[\begin{array}{cc} 2/3 & 1\\ 1 & 2 \end{array} \right]$$

with eigenvalues and eigenvectors:

$$\lambda_1 = 2.53 \quad (\vec{u_1})^t = [0.47, 0.88]$$

 $\lambda_2 = 0.13 \quad (\vec{u_2})^t = [-0.88, 0.47]$

Figure D.3 plots the 3 given points in 2-d space as 'diamonds', as well as their center of gravity (2/3, 1) and the corresponding K-L directions. It uses 'crosses' and 'squares' for the major and minor eigenvector, respectively.



Figure D.3 Example of the Karhunen-Loève transformation, with the 3 points of the Example.

Next we give an implementation of the K-L transform in *Mathematica* [Wol91]. The major step is to compute the covariance matrix C; then, the eigenvalue routine Eigensystem does the rest (and hardest!) of the job.

```
(* given a matrix mat_ with $n$ vectors (rows) of $m$ attributes (columns),
    it creates a matrix with $n$ vectors and their
    first $k$ most 'important' attributes
    (ie., the K-L expansions of these $n$ vectors) *)
KLexpansion[ mat_, k_:2] := mat . Transpose[ KL[mat, k] ];
(* given a matrix with $n$ vectors of $m$ dimensions,
    computes the first $k$ singular vectors,
    ie., the axes of the first $k$ Karhunen-Lo\'{e}ve expansion *)
KL[ mat_ , k_:2 ]:= Module[
    {n,m, avgvec, newmat,i, val, vec },
    {n,m} = Dimensions[mat];
    avgvec = Apply[ Plus, mat] / n //N;
```

```
(* translate vectors, so the mean is zero *)
newmat = Table[ mat[[i]] - avgvec , {i,1,n} ];
{val, vec} = Eigensystem[ Transpose[newmat] . newmat ];
vec[[ Range[1,k] ]]
```

D.2 SVD

As we saw, the eigenvalues and eigenvectors are defined for square matrices. For rectangular matrices, a closely related concept is the *Singular Value Decomposition* (SVD) [Str80, PFTV88, GVL89]: Consider a set of points as before, represented as a $N \times n$ matrix **A**, as in Table D.2.

term	data	information	retrieval	brain	lung
document					
CS-TR1	1	1	1	0	0
CS-TR2	2	2	2	0	0
CS-TR3	1	1	1	0	0
CS-TR4	5	5	5	0	0
MED-TR1	0	0	0	2	2
MED-TR2	0	0	0	3	3
MED-TR3	0	0	0	1	1

Table D.1 Example of a (document-term) matrix

Such a matrix could represent, eg., N patients with n numerical symptoms each (blood pressure, cholesterol level etc), or N sales with n products in a data mining application [AS94], with the dollar amount spent on each product, by the given sale, etc. For concreteness, we shall assume that it represents N documents (rows) with n terms (columns) each, as happens in Information Retrieval (IR) [SFW83], [Dum94]. It would be desirable to group similar documents together, as well as similar terms together. This is exactly what SVD does, automatically! The only 'catch' is that SVD creates a *linear* combination of terms, as opposed to non-linear ones that, eg., Kohonen's neural networks could provide [LSM91, RMS92]. Nevertheless, these groups of terms are valuable: in Information Retrieval terminology, each would correspond to

126

]
a 'concept'; in the Karhunen-Loeve terminology, each group of terms would correspond to an important 'axes'. The formal definition for SVD follows:

Theorem 2 (SVD) Given an $N \times n$ real matrix **A** we can express it as

$$\mathbf{A} = \mathbf{U} \times \mathbf{\Lambda} \times \mathbf{V}^t \tag{D.8}$$

where **U** is a column-orthonormal $N \times r$ matrix, r is the rank of the matrix **A**, **A** is a diagonal $r \times r$ matrix and **V** is a column-orthonormal $k \times r$ matrix.

Proof: See [PTVF92, p. 59].

The entries of Λ are non-negative. If we insist that the diagonal matrix Λ has its elements sorted in descending order, then the decomposition is $unique^1$.

Recall that a matrix **U** is column-orthonormal iff its column vectors are mutually orthogonal and of unit length. Equivalently: $\mathbf{U}^t \times \mathbf{U} = \mathbf{I}$, where **I** is the identity matrix. Schematically, see Figure D.4.



Figure D.4 Illustration of SVD

Eq. D.8 equivalently states that a matrix \mathbf{A} can be brought in the form

$$\mathbf{A} = \lambda_1 \vec{u_1} \times (\vec{v_1})^t + \lambda_2 \vec{u_2} \times (\vec{v_2})^t + \ldots + \lambda_r \vec{u_r} \times (\vec{v_r})^t$$
(D.9)

¹Except when there are equal entries in Λ , in which case they and their corresponding columns of U and V can be permuted.

where $\vec{u_i}$, and $\vec{v_i}$ are column vectors of the **U** and **V** matrices respectively, and λ_i the diagonal elements of the matrix **A**. Intuitively, the SVD identifies 'rectangular blobs' of related values in the **A** matrix. For example, for the above 'toy' matrix of Table D.2, we have two 'blobs' of values, while the rest of the entries are zero. This is confirmed by the SVD, which identifies them both:

Notice that the rank of the matrix is r=2: there are effectively 2 types of documents (CS and Medical documents) and 2 'concepts', i.e., groups-of-terms: the 'CS concept' (that is, the group {'data', 'information', 'retrieval'}), and the 'medical concept' (that is, the group {'lung', 'brain'}). The intuitive meaning of the **U** and **V** matrices is as follows: **U** can be thought of as the *document-to-concept* similarity matrix, while **V**, symmetrically, is the *term-to-concept* similarity matrix. For example, $v_{1,2} = 0$ means that the first term ('data') has zero similarity with the 2nd concept (the 'lung-brain' concept).

The SVD is a powerful operation, with several applications. We list some observations, which are useful for multimedia indexing and Information Retrieval: **Observation D.2** The $N \times N$ matrix $\mathbf{D} = \mathbf{A} \times \mathbf{A}^{t}$ will intuitively give the document-to-document similarities - in our case, it is

	3	6	3	15	0	0	0
	6	12	6	30	0	0	0
	3	6	3	15	0	0	0
$\mathbf{D} = \mathbf{A} \times \mathbf{A}^t = \mathbf{A}$	15	30	15	75	0	0	0
	0	0	0	0	8	12	4
	0	0	0	0	12	18	6
	0	0	0	0	4	6	2

Observation D.3 The eigenvectors of the **D** matrix will be the columns of the **U** matrix of the SVD of \mathbf{A} .

Observation D.4 Symmetrically, the $n \times n$ matrix $\mathbf{T} = \mathbf{A}^t \times \mathbf{A}$ will give the term-to-term similarities - in our example, it is:

	31	31	31	0	0
	31	31	31	0	0
$\mathbf{T} = \mathbf{A}^t \times \mathbf{A}$	31	31	31	0	0
	0	0	0	14	14
	0	0	0	14	14

Observation D.5 Similarly, the eigenvectors of the \mathbf{T} matrix are the columns of the \mathbf{V} matrix of the SVD of \mathbf{A} .

Observation D.6 Both **D** and **T** have the same eigenvalues, who are the squares of the λ_i elements of the Λ matrix of the SVD of **A**.

All the above observations can be proved from Theorem 2 and from the fact that \mathbf{U} and \mathbf{V} are column-orthonormal:

$$\mathbf{A} \times \mathbf{A}^{t} = \mathbf{U} \times \mathbf{\Lambda} \times \mathbf{V}^{t} \times \mathbf{V} \times \mathbf{\Lambda} \times \mathbf{U}^{t} = \mathbf{U} \times \mathbf{\Lambda} \times \mathbf{I} \times \mathbf{\Lambda} \times \mathbf{U}^{t} = \mathbf{U} \times \mathbf{\Lambda}^{2} \times \mathbf{U}^{t} \quad (D.10)$$

and similarly for $\mathbf{A}^t \times \mathbf{A}$. According to (Eq. D.4), the columns of \mathbf{U} are the eigenvectors of $\mathbf{A} \times \mathbf{A}^t$, and its eigenvalues are the diagonal elements of $\mathbf{\Lambda}^2$ (that is, $\lambda_1^2, \dots \lambda_r^2$).

The above observations illustrate the close relationship between the eigenvectors analysis of a matrix, the SVD and the K-L transform, which uses the eigenvectors of the covariance matrix C (Eq. D.7).

It should be noted that the SVD is extremely useful for several settings that involve least-squares optimization, such as in regression, in under-constraint and over-constraint linear problems, etc. See [PTVF92] or [Str80] for more details. Next, we show how it has been applied for Information Retrieval and filtering, under the name of *Latent Semantic Indexing* (LSI).

D.3 SVD AND LSI

Here we discuss SVD in more detail, in the context of Information Filtering. There, SVD has lead to the method of 'Latent Semantic Indexing' (LSI) [FD92b] The idea is to try to group similar terms together, to form a few ($\approx 100 - 300$) 'concepts', and then map the documents into vectors in 'concept'-space, as opposed to vectors in *n*-dimensional space, where *n* is the vocabulary size of the document collection. This approach is a clever, automated way, to take into account term co-occurrences, building effectively a 'thesaurus without semantics': terms that often occur together, are grouped into 'concepts'; every time the user asks for a term, the system determines the relevant 'concepts' and searches for them.

In order to map document or query vectors into to concept space, we need the term-to-concept similarity matrix \mathbf{V} . For example, in the setting of Table D.2, consider the query 'find documents containing the term 'data". In this setting, this query \vec{q} is the vector

$$\vec{q} = \begin{bmatrix} 1\\0\\0\\0\\0 \end{bmatrix}$$
(D.11)

because the first entry corresponds to the term 'data' and the rest to the 4 other terms of our document collection (namely, 'information', 'retrieval', 'brain', 'lung'). To translate \vec{q} to a vector $\vec{q_c}$ in concept space, we need the term-toconcept similarity matrix **V**. The 'translation' is done if we multiply the query vector by \mathbf{V}^t :

$$\vec{q_c} = \mathbf{V}^t \times \vec{q}$$

K-L and SVD

$$= \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$
$$= \begin{bmatrix} 0.58 \\ 0 \end{bmatrix}$$

which correctly corresponds to the fact that the query \vec{q} is rather closely related to the CS group of terms (with 'strength' = 0.58), and unrelated to the medical group of terms ('strength' = 0). What is even more important is that the query $\vec{q_c}$ implicitly involves the terms 'information' and 'retrieval'; thus, an LSI-based system may return documents that do not necessarily contain the requested term 'data', but they are deemed relevant anyway. Eg., according to the running example, the document with the single word 'retrieval' will have the document vector \vec{d}

$$\vec{d} = \begin{bmatrix} 0\\0\\1\\0\\0 \end{bmatrix}$$
(D.12)

which will be mapped to the concept vector $\vec{d_c}$:

$$\vec{d_c} = \mathbf{V}^t \times \vec{d_c} = \begin{bmatrix} 0.58\\0 \end{bmatrix} = \vec{q_c}$$
(D.13)

That is, the above document will be a perfect match for the query ('data'), although it does not contain the query word.

Thanks to its ability to create a thesaurus of co-occurring terms, the LSI method has shown good performance. Experiments in [FD92b] report that LSI has equaled or outperformed standard vector methods and other variants, with improvement of as much as 30% in terms of precision and recall.

D.4 CONCLUSIONS

We have seen some powerful tools, based on eigenvalue analysis. Specifically:

- the 'K-L' transform is the optimal way to do dimensionality reduction: given N vectors with n dimensions, it provides the k most important directions on which to project (k is user defined).
- the SVD (singular value decomposition) operates on an $N \times n$ matrix and groups its rows and columns into r 'similar' groups, sorted in 'strength' order.

Both tools are closely related to the eigenvalue analysis (Eq. D.1): the K-L transform uses the eigenvalues of the covariance matrix; the SVD of a symmetric matrix is identical to its eigenvalue decomposition.

REFERENCES

- [AC75] A.V. Aho and M.J. Corasick. Fast pattern matching: an aid to bibliographic search. *CACM*, 18(6):333-340, June 1975.
- [ACF+93] Manish Arya, William Cody, Christos Faloutsos, Joel Richardson, and Arthur Toga. Qbism: a prototype 3-d medical image database system. IEEE Data Engineering Bulletin, 16(1):38-42, March 1993.
- [ACF+94] Manish Arya, William Cody, Christos Faloutsos, Joel Richardson, and Arthur Toga. Qbism: Extending a dbms to support 3d medical images. Tenth Int. Conf. on Data Engineering (ICDE), pages 314– 325, February 1994.
- [AFS93] Rakesh Agrawal, Christos Faloutsos, and Arun Swami. Efficient similarity search in sequence databases. In Foundations of Data Organization and Algorithms (FODO) Conference, Evanston, Illinois, October 1993. also available through anonymous ftp, from olympos.cs.umd.edu: ftp/pub/TechReports/fodo.ps.
- [AGI+92] Rakesh Agrawal, Sakti Ghosh, Tomasz Imielinski, Bala Iyer, and Arun Swami. An interval classifier for database mining applications. *VLDB Conf. Proc.*, pages 560-573, August 1992.
- [AGM⁺90] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, and D.J. Lipman. A basic local alignment search tool. Journal of Molecular Biology, 215(3):403-410, 1990.
- [AIS93a] Rakesh Agrawal, Tomasz Imielinski, and Arun Swami. Database mining: a performance perspective. IEEE Trans. on Knowledge and Data Engineering, 5(6):914-925, 1993.
- [AIS93b] Rakesh Agrawal, Tomasz Imielinski, and Arun Swami. Mining association rules between sets of items in large databases. Proc. ACM SIGMOD, pages 207-216, May 1993.
- [AKW88] Alfred V. Aho, Brian W. Kernighan, and Peter J. Weinberger. The AWK Programming Language. Addison Wesley, 1988.

- [AS94] Rakesh Agrawal and Ramakrishnan Srikant. Fast algorithms for mining association rules in large databases. Proc. of VLDB Conf., pages 487-499, September 1994.
- [BB82] D. Ballard and C. Brown. Computer Vision. Prentice Hall, 1982.
- [BCC94] Eric W. Brown, James P. Callan, and W. Bruce Croft. Fast incremental indexing for full-text information retrieval. Proc. of VLDB Conf., pages 192-202, September 1994.
- [Ben75] J.L. Bentley. Multidimensional Binary Search Trees Used for Associative Searching. CACM, 18(9):509-517, September 1975.
- [Ben79] Jon L. Bentley. Multidimensional binary search trees in database applications. *IEEE Trans. on Software Engineering (TSE)*, SE-5(4):333-340, July 1979.
- [BF95] Alberto Belussi and Christos Faloutsos. Estimating the selectivity of spatial queries using the 'correlation' fractal dimension. Proc. of VLDB, pages 299-310, September 1995.
- [Bia69] T. Bially. Space-filling curves: Their generation and their application to bandwidth reduction. *IEEE Trans. on Information Theory*, IT-15(6):658-664, November 1969.
- [BJR94] George E.P. Box, Gwilym M. Jenkins, and Gregory C. Reinsel. Time Series Analysis: Forecasting and Control. Prentice Hall, Englewood Cliffs, NJ, 1994. 3rd Edition.
- [BKS93] Thomas Brinkhoff, Hans-Peter Kriegel, and Bernhard Seeger. Efficient processing of spatial joins using r-trees. Proc. of ACM SIG-MOD, pages 237-246, May 1993.
- [BKSS90] N. Beckmann, H.-P. Kriegel, R. Schneider, and B. Seeger. The r*-tree: an efficient and robust access method for points and rectangles. ACM SIGMOD, pages 322-331, May 1990.
- [BKSS94] Thomas Brinkhoff, Hans-Peter Kriegel, Ralf Schneider, and Bernhard Seeger. Multi-step processing of spatial joins. ACM SIGMOD, pages 197-208, May 1994.
- [BM72] R. Bayer and E. McCreight. Organization and maintenance of large ordered indexes. Acta Informatica, 1(3):173-189, 1972.
- [BM77] R.S. Boyer and J.S. Moore. A fast string searching algorithm. CACM, 20(10):762-772, October 1977.

- [But71] Arthur R. Butz. Alternative algorithm for hilbert's space-filling curve. *IEEE Trans. on Computers*, C-20(4):424-426, April 1971.
- [BYG92] Ricardo Baeza-Yates and Gaston H. Gonnet. A new approach to text searching. Comm. of ACM (CACM), 35(10):74-82, October 1992.
- [CE92] M. Castagli and S. Eubank. Nonlinear Modeling and Forecasting. Addison Wesley, 1992. Proc. Vol. XII.
- [CF84] S. Christodoulakis and C. Faloutsos. Design considerations for a message file server. *IEEE Trans. on Software Engineering*, SE-10(2):201-210, March 1984.
- [Cod70] E. F. Codd. A relational model of data for large shared data banks. Comm. of ACM, 13(6):377-387, 1970.
- [Coo70] W.S. Cooper. On deriving design equations for information retrieval systems. JASIS, pages 385–395, November 1970.
- [CoPES92] Mathematical Committee on Physical and NSF Engineering Sciences. Grand Challenges: High Performance Computing and Communications. National Science Foundation, 1992. The FY 1992 U.S. Research and Development Program.
- [CP90] Doug Cutting and Jan Pedersen. Optimizations for dynamic inverted index maintenance. *Proc. SIGIR*, pages 405-411, 1990.
- [Cra94] Richard E. Crandall. Projects in Scientific Computation. Springer-Verlag New York, Inc., 1994.
- [CS89] W.W. Chang and H.J. Schek. A signature access method for the starbust database system. In Proc. VLDB Conference, pages 145– 153, Amsterdam, Netherlands, August 1989.
- [CSY87] Shi-Kuo Chang, Qing-Yun Shi, and Cheng-Wen Yan. Iconic Indexing by 2-D Strings. IEEE Transactions on Pattern Analysis and Machine Intelligence, 9(3):413-428, May 1987.
- [CTH⁺86] S. Christodoulakis, M. Theodoridou, F. Ho, M. Papa, and A. Pathria. Multimedia document presentation, information extraction and document formation in minos: a model and a system. ACM TOOIS, 4(4), October 1986.
- [Dat86] Chris J. Date. An Introduction to Database Systems. Addison-Wesley, 1986. Vol. I; 4th ed.

- [Dau92] Ingrid Daubechies. Ten Lectures on Wavelets. Capital City Press, Montpelier, Vermont, 1992. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA.
- [DDF⁺90] S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman. Indexing by latent semantic analysis. *Journal of the American Society for Information Science*, 41(6):391-407, September 1990.
- [DH73a] R. Duda and P Hart. Pattern Classification and Scene Analysis. Wiley, New York, 1973.
- [DH73b] R.O. Duda and P.E. Hart. Pattern Classification and Scene Analysis. Wiley, New York, 1973.
- [Dum94] Susan T. Dumais. Latent semantic indexing (lsi) and trec-2. In D. K. Harman, editor, The Second Text Retrieval Conference (TREC-2), pages 105-115, Gaithersburg, MD, March 1994. NIST. Special publication 500-215.
- [Dye82] C.R. Dyer. The space efficiency of quadtrees. Computer Graphics and Image Processing, 19(4):335-348, August 1982.
- [Eli75] P. Elias. Universal codeword sets and representations of integers. IEEE Trans. on Information Theory, IT-21:194-203, 1975.
- [EM66] Robert D. Edwards and John Magee. Technical Analysis of Stock Trends. John Magee, Springfield, Massachusetts, 1966. 5th Edition, second printing.
- [EMS⁺86] J. Ewing, S. Mehrabanzad, S. Sheck, D. Ostroff, and B. Shneiderman. An experimental comparison of a mouse and arrow-jump keys for an interactive encyclopedia. *Int. Journal of Man-Machine Studies*, 24(1):29-45, January 1986.
- [Equ93] W. Equitz. Retrieving images from a database using texture algorithms from the QBIC system. Research report, IBM Almaden Research Center, San Jose, CA, 1993.
- [Fal85] C. Faloutsos. Access methods for text. ACM Computing Surveys, 17(1):49-74, March 1985.
- [Fal88] C. Faloutsos. Gray codes for partial match and range queries. IEEE Trans. on Software Engineering, 14(10):1381-1393, October 1988. early version available as UMIACS-TR-87-4, also CS-TR-1796.

- [Fal92a] C. Faloutsos. Analytical results on the quadtree decomposition of arbitrary rectangles. Pattern Recognition Letters, 13(1):31-40, January 1992.
- [Fal92b] Christos Faloutsos. Signature files. In William Bruce Frakes and Ricardo Baeza-Yates, editors, Information Retrieval: Data Structures and Algorithms. Prentice Hall, Englewood Cliffs, NJ, 1992.
- [FBF⁺94] C. Faloutsos, R. Barber, M. Flickner, J. Hafner, W. Niblack, D. Petkovic, and W. Equitz. Efficient and effective querying by image content. *Journal of Intell. Inf. Systems*, 3(3/4):231-262, July 1994.
- [FBY92] W. Frakes and R. Baeza-Yates. Information Retrieval: Data Structures and Algorithms. Prentice-Hall, 1992.
- [FC87] C. Faloutsos and S. Christodoulakis. Optimal signature extraction and information loss. ACM TODS, 12(3):395-428, September 1987.
- [FD92a] Peter W. Foltz and Susan T. Dumais. Personalized information delivery: An analysis of information filtering methods. Communications of the ACM, 35(12):51-60, December 1992.
- [FD92b] Peter W. Foltz and Susan T. Dumais. Personalized information delivery: an analysis of information filtering methods. Comm. of A CM (CA CM), 35(12):51-60, December 1992.
- [FG96] Christos Faloutsos and Volker Gaede. Analysis of the z-ordering method using the hausdorff fractal dimension. VLDB, September 1996.
- [FH69] J.R. Files and H.D. Huskey. An information retrieval system based on superimposed coding. Proc. AFIPS FJCC, 35:423-432, 1969.
- [Fie93] D.J. Field. Scale-invariance and self-similar 'wavelet' transforms: an analysis fo natural scenes and mammalian visual systems. In M. Farge, J.C.R. Hunt, and J.C. Vassilicos, editors, Wavelets, Fractals, and Fourier Transforms, pages 151-193. Clarendon Press, Oxford, 1993.
- [FJM94] Christos Faloutsos, H.V. Jagadish, and Yannis Manolopoulos. Analysis of the n-dimensional quadtree decomposition for arbitrary hyper-rectangles. CS-TR-3381, UMIACS-TR-94-130, Dept. of Computer Science, Univ. of Maryland, College Park, MD, December 1994. to appear in IEEE TKDE.

- [FK94] Christos Faloutsos and Ibrahim Kamel. Beyond Uniformity and Irndependence: Analysis of R-trees Using the Concept of Fractal Dimension. Proc. ACM SIGACT-SIGMOD-SIGART PODS, pages 4-13, May 1994. Also available as CS-TR-3198, UMIACS-TR-93-130.
- [FL95] Christos Faloutsos and King-Ip (David) Lin. Fastmap: a fast algorithm for indexing, data-mining and visualization of traditional and multimedia datasets. Proc. of ACM-SIGMOD, pages 163-174, May 1995.
- [FN75] Keinosuke Fukunaga and Patrenahalli M. Narendra. A branch and bound algorithm for computing k-nearest neighbors. *IEEE Trans.* on Computers (TOC), C-24(7):750-753, July 1975.
- [FNPS79] R. Fagin, J. Nievergelt, N. Pippenger, and H.R. Strong. Extendible hashing - a fast access method for dynamic files. ACM TODS, 4(3):315-344, September 1979.
- [FR89a] C. Faloutsos and W. Rego. Tri-cell: a data structure for spatial objects. *Information Systems*, 14(2):131-139, 1989. early version available as UMIACS-TR-87-15, CS-TR-1829.
- [FR89b] C. Faloutsos and S. Roseman. Fractals for secondary key retrieval. Eighth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS), pages 247-252, March 1989. also available as UMIACS-TR-89-47 and CS-TR-2242.
- [FR91] C. Faloutsos and Y. Rong. Dot: a spatial access method using fractals. In *IEEE Data Engineering Conference*, pages 152–159, Kobe, Japan, April 1991. early version available as UMIACS-TR-89-31, CS-TR-2214.
- [Fre60] E. Fredkin. Trie memory. CACM, 3(9):490-500, September 1960.
- [FRM94] Christos Faloutsos, M. Ranganathan, and Yannis Manolopoulos. Fast subsequence matching in time-series databases. Proc. ACM SIGMOD, pages 419–429, May 1994. 'Best Paper' award; also available as CS-TR-3190, UMIACS-TR-93-131, ISR TR-93-86.
- [FSN+95] Myron Flickner, Harpreet Sawhney, Wayne Niblack, Jon Ashley, Qian Huang, Byron Dom, Monika Gorkani, Jim Hafner, Denis Lee, Dragutin Petkovic, David Steele, and Peter Yanker. Query by image and video content: the qbic system. *IEEE Computer*, 28(9):23-32, September 1995.

- [FSR87] C. Faloutsos, T. Sellis, and N. Roussopoulos. Analysis of object oriented spatial access methods. Proc. ACM SIGMOD, pages 426– 439 426–439, May 1987. also available as SRC-TR-87-30, UMIACS-TR-86-27, CS-TR-1781.
- [Fuk90] Keinosuke Fukunaga. Introduction to Statistical Pattern Recognition. Academic Press, 1990. 2nd Edition.
- [Gae95] V. Gaede. Optimal redundancy in spatial database systems. In Proc. 4th Int. Symp. on Spatial Databases (SSD'95), pages 96-116, 1995.
- [Gal91] D. Le Gall. Mpeg: a video compression standard for multimedia applications. Comm. of ACM (CACM), 34(4):46-58, April 1991.
- [Gar82] I. Gargantini. An effective way to represent quadtrees. Comm. of ACM (CACM), 25(12):905-910, December 1982.
- [GBYS92] Gaston H. Gonnet, Ricardo A. Baeza-Yates, and Tim Snider. New indices for text: Pat trees and pat arrays. In William B. Frakes and Ricardo Baeza-Yates, editors, *Information Retrieval*. Prentice Hall, 1992.
- [GG95] Volker Gaede and Oliver Guenther. Survey on multidimensional access methods. Technical Report ISS-16, Institut fuer Wirtschaftsinformatik, Humboldt-Universitaet zu Berlin, August 1995.
- [GGMT94] Louis Gravano, Hector Garcia-Molina, and Anthony Tomasic. The effectiveness of gloss for the text database discovery problem. ACM SIGMOD, pages 126-137, May 1994.
- [GK95] Dina Q. Goldin and Paris C. Kanellakis. On similarity queries for time-series data: Constraint specification and implementation. Int. Conf. on Principles and Practice of Constraint Programming (CP95), September 1995.
- [GR94] V. Gaede and W.-F. Riekert. Spatial access methods and query processing in the object-oriented GIS GODOT. In Proc. of the AGDM'94 Workshop, pages 40-52, Delft, The Netherlands, 1994. Netherlands Geodetic Commission.
- [GR95] Venkat N. Gudivada and Vijay V. Raghavan. Content-based image retrieval systems. IEEE Computer, 28(9):18-22, September 1995.

- [GT87] G.H. Gonnet and F.W. Tompa. Mind your grammar: a new approach to modelling text. Proc. of the Thirteenth Int. Conf. on Very Large Data Bases, pages 339-346, September 1987.
- [Gun86] O. Gunther. The cell tree: an index for geometric data. Memorandum No. UCB/ERL M86/89, Univ. of California, Berkeley, December 1986.
- [Gut84] A. Guttman. R-trees: a dynamic index structure for spatial searching. Proc. ACM SIGMOD, pages 47-57, June 1984.
- [GVL89] G. H. Golub and C. F. Van-Loan. Matrix Computations. The Johns Hopkins University Press, Baltimore, 1989. 2nd edition.
- [Har71] M.C. Harrison. Implementation of the substring test by hashing. CACM, 14(12):777-779, December 1971.
- [Har94] D. Harman. The second text retrieval conference (trec-2). Special Publication 500-215, National Institute of Standards and Technology, Gaithersburg, MD., 1994.
- [Has81] R.L. Haskin. Special-purpose processors for text retrieval. Database Engineering, 4(1):16-29, September 1981.
- [HD80] P.A.V. Hall and G.R. Dowling. Approximate string matching. ACM Computing Surveys, 12(4):381-402, December 1980.
- [HH83] R.L. Haskin and L.A. Hollaar. Operational characteristics of a hardware-based pattern matcher. ACM TODS, 8(1):15-40, March 1983.
- [HN83] K. Hinrichs and J. Nievergelt. The grid file: a data structure to support proximity queries on spatial objects. Proc. of the WG'83 (Intern. Workshop on Graph Theoretic Concepts in Computer Science), pages 100-113, 1983.
- [Hol79] L.A. Hollaar. Text retrieval computers. IEEE Computer Magazine, 12(3):40-50, March 1979.
- [Hor86] Berthold Horn. Robot Vision. MIT Press, Cambridge, Mass., 1986.
- [HS79] G.M. Hunter and K. Steiglitz. Operations on images using quad trees. IEEE Trans. on PAMI, PAMI-1(2):145-153, April 1979.
- [HS91] Andrew Hume and Daniel Sunday. Fast string searching. Software - Practice and Experience, 21(11):1221-1248, November 1991.

- [HS95] M. Houtsma and A. Swami. Set-oriented mining for association rules. In Proceedings of IEEE Data Engineering Conference, March 1995. Also appeared as IBM Research Report RJ 9567.
- [HSW88] A. Hutflesz, H.-W. Six, and P. Widmayer. Twin grid files: Space optimizing access schemes. Proc. of ACM SIGMOD, pages 183-190, June 1988.
- [HU79] J.E. Hopcroft and J.D. Ullman. Introduction to Automata Theory, Languages, and Computation. Addison Wesley, Reading, Mass., 1979.
- [Jag90a] H.V. Jagadish. Linear clustering of objects with multiple attributes. A CM SIGMOD Conf., pages 332-342, May 1990.
- [Jag90b] H.V. Jagadish. Spatial search with polyhedra. Proc. Sixth IEEE Int'l Conf. on Data Engineering, February 1990.
- [Jag91] H.V. Jagadish. A retrieval technique for similar shapes. Proc. ACM SIGMOD Conf., pages 208-217, May 1991.
- [JMM95] H.V. Jagadish, Alberto O. Mendelzon, and Tova Milo. Similaritybased queries. Proc. ACM SIGACT-SIGMOD-SIGART PODS, pages 36-45, May 1995.
- [JN92] Ramesh Jain and Wayne Niblack. NSF Workshop on Visual Information Management, February 1992.
- [JS89] Theodore Johnson and Dennis Shasha. Utilization of b-trees with inserts, deletes and modifies. *Proc. of ACM SIGACT-SIGMOD-SIGART PODS*, pages 235-246, March 1989.
- [Jur92] Ronald K. Jurgen. Digital video. IEEE Spectrum, 29(3):24-30, March 1992.
- [KF93] Ibrahim Kamel and Christos Faloutsos. On packing r-trees. Second Int. Conf. on Information and Knowledge Management (CIKM), November 1993.
- [KF94] Ibrahim Kamel and Christos Faloutsos. Hilbert R-tree: An Improved R-tree Using Fractals. In Proceedings of VLDB Conference,, pages 500-509, Santiago, Chile, September 1994.
- [Kim88] R.E. Kimbrell. Searching for text? send an n-gram! Byte, 13(5):297-312, May 1988.

- [KMP77] D.E. Knuth, J.H. Morris, and V.R. Pratt. Fast pattern matching in strings. SIAM J. Comput, 6(2):323-350, June 1977.
- [Kno75] G.D. Knott. Hashing functions. Computer Journal, 18(3):265-278, 1975.
- [Knu73] D.E. Knuth. The Art of Computer Programming, Vol. 3: Sorting and Searching. Addison-Wesley, Reading, Mass, 1973.
- [KR87] R.M. Karp and M.O. Rabin. Efficient randomized pattern-matching algorithms. IBM Journal of Research and Development, 31(2):249-260, March 1987.
- [KS91] Henry F. Korth and Abraham Silberschatz. Database System Concepts. McGraw Hill, 1991.
- [KSF⁺96] Flip Korn, Nikolaos Sidiropoulos, Christos Faloutsos, Eliot Siegel, and Zenon Protopapas. Fast nearest-neighbor search in medical image databases. Conf. on Very Large Data Bases (VLDB), September 1996. Also available as Univ. of Maryland tech. report: CS-TR-3613, ISR-TR-96-13.
- [KTF95] Anil Kumar, Vassilis J. Tsotras, and Christos Faloutsos. Access methods for bi-temporal databases. Int. Workshop on Temporal Databases, September 1995.
- [KW78] Joseph B. Kruskal and Myron Wish. Multidimensional scaling. SAGE publications, Beverly Hills, 1978.
- [Lar78] P. Larson. Dynamic hashing. BIT, 18:184–201, 1978.
- [Lar82] P.A. Larson. Performance analysis of linear hashing with partial expansions. ACM TODS, 7(4):566-587, December 1982.
- [Lar85] P.A. Larson. Hash files: Some recent developments. In Proc. of the First Intern. Conference on Supercomputing Systems, pages 671-679, St. Petersburg, Florida, December 1985.
- [Lar88] P.-A. Larson. Dynamic hash tables. Comm. of ACM (CACM), 31(4):446-457, April 1988.
- [LeB92] Blake LeBaron. Nonlinear forecasts for the s\&p stock index. In M. Castagli and S. Eubank, editors, Nonlinear Modeling and Forecasting, pages 381-393. Addison Wesley, 1992. Proc. Vol. XII.
- [Les78] M.E. Lesk. Some Applications of Inverted Indexes on the UNIX System. Bell Laboratories, Murray Hill, New Jersey, 1978.

- [Lit80] W. Litwin. Linear hashing: a new tool for file and table addressing. In Proc. 6th International Conference on VLDB, pages 212-223, Montreal, October 1980.
- [LL89] D.L. Lee and C.-W. Leng. Partitioned signature file: Designs and performance evaluation. ACM Trans. on Information Systems (TOIS), 7(2):158-180, April 1989.
- [LR94] Ming-Ling Lo and Chinya V. Ravishankar. Spatial joins using seeded trees. ACM SIGMOD, pages 209-220, May 1994.
- [LS90] David B. Lomet and Betty Salzberg. The hb-tree: a multiattribute indexing method with good guaranteed performance. ACM TODS, 15(4):625-658, December 1990.
- [LSM91] Xia Lin, Dagobert Soergel, and Gary Marchionini. A self-organizing semantic map for information retrieval. Proc. of ACM SIGIR, pages 262-269, October 1991.
- [Lum70] V.Y. Lum. Multi-attribute retrieval with combined indexes. CACM, 13(11):660-665, November 1970.
- [LW75] R. Lowerance and R.A. Wagner. An extension of the string-tostring correction problem. JACM, 22(2):3-14, April 1975.
- [LW89] Carl E. Langenhop and William E. Wright. A model of the dynamic behavior of b-trees. *Acta Informatica*, 27:41–59, 1989.
- [Man77] B. Mandelbrot. Fractal Geometry of Nature. W.H. Freeman, New York, 1977.
- [Mar79] G.N.N. Martin. Spiral storage: Incrementally augmentable hash addressed storage. Theory of Computation, Report No. 27, Univ. of Warwick, Coventry, England, March 1979.
- [McI82] M.D. McIlroy. Development of a spelling list. IEEE Trans. on Communications, COM-30(1):91-99, January 1982.
- [MJFS96] Bongki Moon, H.V. Jagadish, Christos Faloutsos, and Joel H. Saltz. Analysis of the clustering properties of hilbert space-filling curve. Technical Report CS-TR-3611, Dept. of Computer Science, Univ. of Maryland, College Park, MD, 1996.
- [ML86] Lothar M. Mackert and Guy M. Lohman. R* optimizer validation and performance evaluation for distributed queries. *Proc. of 12th Int. Conf. on Very Large Data Bases (VLDB)*, August 1986.

[Moo49] C. Mooers. Application of random codes to the gathering of statistical information. Bulletin 31, Zator Co, Cambridge, Mass, 1949. based on M.S. thesis, MIT, January 1948.

144

- [Mor68] Donald R. Morrison. Patricia practical algorithm to retrieve information coded in alphanumeric. Journal of ACM (JACM), 15(4):514-534, October 1968.
- [MRT91] Carlo Meghini, Fausto Rabitti, and Constantino Thanos. Conceptual modeling of multimedia documents. *IEEE Computer*, 24(10):23-30, October 1991.
- [Mur83] F. Murtagh. A survey of recent advances in hierarchical clustering algorithms. *The Computer Journal*, 26(4):354-359, 1983.
- [MW94] Udi Manber and Sun Wu. Glimpse: a tool to search through entire file systems. Proc. of USENIX Techn. Conf., 1994. Also available as TR 93-94, Dept. of Comp. Sc., Univ. of Arizona, Tucson, or through anonymous ftp (ftp://cs.arizona.edu/glimpse/glimpse.ps.Z).
- [NBE⁺93] Wayne Niblack, Ron Barber, Will Equitz, Myron Flickner, Eduardo Glasman, Dragutin Petkovic, Peter Yanker, Christos Faloutsos, and Gabriel Taubin. The qbic project: Querying images by content using color, texture and shape. SPIE 1993 Intl. Symposium on Electronic Imaging: Science and Technology, Conf. 1908, Storage and Retrieval for Image and Video Databases, February 1993. Also available as IBM Research Report RJ 9203 (81511), Feb. 1, 1993, Computer Science.
- [NC91] A. Desai Narasimhalu and Stavros Christodoulakis. Multimedia information systems: the unfolding of a reality. *IEEE Computer*, 24(10):6-8, October 1991.
- [NHS84] J. Nievergelt, H. Hinterberger, and K.C. Sevcik. The grid file: an adaptable, symmetric multikey file structure. ACM TODS, 9(1):38-71, March 1984.
- [Nof86] P.J. Nofel. 40 million hits on optical disk. Modern Office Technology, pages 84-88, March 1986.
- [ODL93] Katia Obraczka, Peter B. Danzig, and Shih-Hao Li. Internet resourse discovery services. *IEEE Computer*, September 1993.
- [OM84] J.A. Orenstein and T.H. Merrett. A class of data structures for associative searching. Proc. of SIGACT-SIGMOD, pages 181-190, April 1984.

- [OM88] J.A. Orenstein and F.A. Manola. Probe spatial data modeling and query processing in an image database application. *IEEE Trans.* on Software Engineering, 14(5):611-629, May 1988.
- [Ore86] J. Orenstein. Spatial query processing in an object-oriented database system. *Proc. ACM SIGMOD*, pages 326-336, May 1986.
- [Ore89] J.A. Orenstein. Redundancy in spatial databases. Proc. of ACM SIGMOD Conf., May 1989.
- [Ore90] J.A. Orenstein. A comparison of spatial query processing techniques for native and parameter spaces. Proc. of ACM SIGMOD Conf., pages 343-352, 1990.
- [OS75] Alan Victor Oppenheim and Ronald W. Schafer. Digital Signal Processing. Prentice-Hall, Englewood Cliffs, N.J., 1975.
- [OS95] Virginia E. Ogle and Michael Stonebraker. Chabot: Retrieval from a relational database of images. *IEEE Computer*, 28(9):40-48, September 1995.
- [Pet80] J.L. Peterson. Computer programs for detecting and correcting spelling errors. CACM, 23(12):676-687, December 1980.
- [PF94] Euripides G.M. Petrakis and Christos Faloutsos. Similarity searching in large image databases, 1994. submitted for publication. Also available as technical report at MUSIC with # TR-01-94.
- [PF96] Euripides G.M. Petrakis and Christos Faloutsos. Similarity searching in medical image databases. *IEEE Trans. on Knowledge and Data Engineering (TDKE)*, 1996. To appear. Also available as technical report at MUSIC with # TR-01-94, UMIACS-TR-94-134, CS-TR-3388.
- [PFTV88] William H. Press, Brian P. Flannery, Saul A. Teukolsky, and William T. Vetterling. Numerical Recipes in C. Cambridge University Press, 1988.
- [PO95] E. G.M. Petrakis and S. C. Orphanoudakis. A Generalized Approach for Image Indexing and Retrieval Based Upon 2-D Strings. In S.-K. Chang, E. Jungert, and G. Tortora, editors, Intelligent Image Database Systems - Spatial Reasoning, Image Indexing and Retrieval using Symbolic Projections. World Scientific Pub. Co., 1995. To be publised. Also available as FORTH-ICS/TR-103.

- [Pri84] Joseph Price. The optical disk pilot project at the library of congress. Videodisc and Optical Disk, 4(6):424-432, November 1984.
- [PSTW93] B. Pagel, H. Six, H. Toben, and P. Widmayer. Towards an analysis of range query performance. Proc. of ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS), pages 214-221, May 1993.
- [PTVF92] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. Numerical Recipes in C. Cambridge University Press, 1992. 2nd Edition.
- [Ras92] Edie Rasmussen. Clustering algorithms. In William B. Frakes and Ricardo Baeza-Yates, editors, Information Retrieval: Data Structures and Algorithms, pages 419-442. Prentice Hall, 1992.
- [RBC+92] Mary Beth Ruskai, Gregory Beylkin, Ronald Coifman, Ingrid Daubechies, Stephane Mallat, Yves Meyer, and Louise Raphael. Wavelets and Their Applications. Jones and Bartlett Publishers, Boston, MA, 1992.
- [RF91] Yi Rong and Christos Faloutsos. Analysis of the clustering property of peano curves. Techn. Report CS-TR-2792, UMIACS-TR-91-151, Univ. of Maryland, December 1991.
- [Riv76] R.L. Rivest. Partial match retrieval algorithms. SIAM J. Comput, 5(1):19-50, March 1976.
- [RJ93] Lawrence Rabiner and Biing-Hwang Juang. Fundamentals of Speech Recognition. Prentice Hall, 1993.
- [RKV95] Nick Roussopoulos, Steve Kelley, and F. Vincent. Nearest Neighbor Queries. Proc. of ACM-SIGMOD, pages 71-79, May 1995.
- [RL85] N. Roussopoulos and D. Leifker. Direct spatial search on pictorial databases using packed R-trees. Proc. ACM SIGMOD, May 1985.
- [RMS92] Helge Ritter, Thomas Martinetz, and Klaus Schulten. Neural Computation and Self-Organizing Maps. Addison Wesley, Reading, MA, 1992.
- [Rob79] C.S. Roberts. Partial-match retrieval via the method of superimposed codes. Proc. IEEE, 67(12):1624-1642, December 1979.

- [Rob81] J.T. Robinson. The k-d-b-tree: a search structure for large multidimensional dynamic indexes. Proc. ACM SIGMOD, pages 10-18, 1981.
- [Roc71] J.J. Rocchio. Performance indices for document retrieval. In G. Salton, editor, The SMART Retrieval System - Experiments in Automatic Document Processing. Prentice-Hall Inc, Englewood Cliffs, New Jersey, 1971. Chapter 3.
- [RS92] Fausto Rabitti and Pascuale Savino. An information retrieval approach for image databases. In VLDB Conf. Proceedings, pages 574–584, Vancouver, BC, Canada, August 1992.
- [RV91] Oliver Rioul and Martin Vetterli. Wavelets and signal processing. IEEE SP Magazine, pages 14-38, October 1991.
- [Sal71a] G. Salton. Relevance feedback and the optimization of retrieval effectiveness. In G. Salton, editor, The SMART Retrieval System -Experiments in Automatic Document Processing. Prentice-Hall Inc, Englewood Cliffs, New Jersey, 1971. Chapter 15.
- [Sal71b] G. Salton. The SMART Retrieval System Experiments in Automatic Document Processing. Prentice-Hall Inc, Englewood Cliffs, New Jersey, 1971.
- [Sch91] Manfred Schroeder. Fractals, Chaos, Power Laws: Minutes From an Infinite Paradise. W.H. Freeman and Company, New York, 1991.
- [SD76] D.G. Severance and R.A. Duhne. A practitioner's guide to addressing algorithms. CACM, 19(6):314-326, 1976.
- [SDKR87] R. Sacks-Davis, A. Kent, and K. Ramamohanarao. Multikey access methods based on superimposed coding techniques. ACM Trans. on Database Systems (TODS), 12(4):655-696, December 1987.
- [SDR83] R. Sacks-Davis and K. Ramamohanarao. A two level superimposed coding scheme for partial match retrieval. Information Systems, 8(4):273-280, 1983.
- [SFW83] G. Salton, E.A. Fox, and H. Wu. Extended boolean information retrieval. CACM, 26(11):1022-1036, November 1983.
- [Sha88] C.A. Shaffer. A formula for computing the number of quadtree node fragments created by a shift. Pattern Recognition Letters, 7(1):45-49, January 1988.

- [Sig93] Karl Sigmund. Games of Life: Explorations in Ecology, Evolution and Behaviour. Oxford University Press, 1993.
- [SK83] David Sankoff and Joseph B. Kruskal. Time Warps, String Edits and Macromolecules: the Theory and Practice of Sequence Comparisons. Addison-Wesley Publishing Company, Inc., Reading, MA, 1983.
- [SK86] C. Stanfill and B. Kahle. Parallel free-text search on the connection machine system. *CACM*, 29(12):1229–1239, December 1986.
- [SK90] Bernhard Seeger and Hans-Peter Kriegel. The buddy-tree: an efficient and robust access method for spatial database systems. *Proc.* of *VLDB*, pages 590-601, August 1990.
- [SL76] D.G. Severance and G.M. Lohman. Differential files: Their application to the maintenance of large databases. ACM TODS, 1(3):256-267, September 1976.
- [SM83] G. Salton and M.J. McGill. Introduction to Modern Information Retrieval. McGraw-Hill, 1983.
- [SOF+92] G.A. Story, L. O'Gorman, D.S. Fox, L.L. Schaper, and H.V. Jagadish. The rightpages: an electronic library for alerting and browsing. *IEEE Computer*, 25(9):17-26, September 1992.
- [SS88] R. Stam and Richard Snodgrass. A bibliography on temporal databases. IEEE Bulletin on Data Engineering, 11(4), December 1988.
- [SSH86] M. Stonebraker, T. Sellis, and E. Hanson. Rule indexing implementations in database systems. In Proceedings of the First International Conference on Expert Database Systems, Charleston, SC, April 1986.
- [SSN87] C.A. Shaffer, H. Samet, and R.C. Nelson. Quilt: a geographic information system based on quadtrees. Technical Report CS-TR-1885.1, Univ. of Maryland, Dept. of Computer Science, July 1987. to appear in the International Journal of Geographic Information Systems.
- [ST84] J. Stubbs and F.W. Tompa. Waterloo and the new oxford english dictionary project. In Proc. of the Twentieth Annual Conference on Editorial Problems, Toronto, Ontario, November 1984. in press.
- [Sta80] T.A. Standish. Data Structure Techniques. Addison Wesley, 1980.

- [Sti60] S. Stiassny. Mathematical analysis of various superimposed coding methods. American Documentation, 11(2):155-169, February 1960.
- [Str80] Gilbert Strang. Linear Algebra and its Applications. Academic Press, 1980. 2nd edition.
- [Sun90] D.M. Sunday. A very fast substring search algorithm. Comm. of A CM (CA CM), 33(8):132-142, August 1990.
- [SW78] G. Salton and A. Wong. Generation and search of clustered files. ACM TODS, 3(4):321-346, December 1978.
- [TC83] D. Tsichritzis and S. Christodoulakis. Message files. ACM Trans. on Office Information Systems, 1(1):88-98, January 1983.
- [TGMS94] Anthony Tomasic, Hector Garcia-Molina, and Kurt Shoens. Incremental updates of inverted lists for text document retrieval. ACM SIGMOD, pages 289-300, May 1994.
- [Ton90] Howell Tong. Non-Linear Time Series: a Dynamical System Approach. Clarendon Press, Oxford, 1990.
- [TSW⁺85] G.R. Thoma, S. Suthasinekul, F.A. Walker, J. Cookson, and M. Rashidian. A prototype system for the electronic storage and retrieval of document images. ACM TOOIS, 3(3), July 1985.
- [Vas93] Dimitris Vassiliadis. The input-state space approach to the prediction of auroral geomagnetic activity from solar wind variables. Int. Workshop on Applications of Artificial Intelligence in Solar Terrestrial Physics, September 1993.
- [VM] Brani Vidakovic and Peter Mueller. Wavelets for Kids. Duke University, Durham, NC. ftp://ftp.isds.duke.edu/pub/Users/brani/papers/.
- [VR79] C.J. Van-Rijsbergen. Information Retrieval. Butterworths, London, England, 1979. 2nd edition.
- [Wal91] Gregory K. Wallace. The jpeg still picture compression standard. CACM, 34(4):31-44, April 1991.
- [WG94] Andreas S. Weigend and Neil A. Gerschenfeld. Time Series Prediction: Forecasting the Future and Understanding the Past. Addison Wesley, 1994.

- [Whi81] M. White. N-Trees: Large Ordered Indexes for Multi-Dimensional Space. Application Mathematics Research Staff, Statistical Research Division, U.S. Bureau of the Census, December 1981.
- [WM92] Sun Wu and Udi Manber. Text searching allowing errors. Comm. of ACM (CACM), 35(10):83-91, October 1992.
- [Wol91] Stephen Wolfram. *Mathematica*. Addison Wesley, 1991. Second Edition.
- [WS93] Kuansan Wang and Shihab Shamma. Spectral shape analysis in the central auditory system. *NNSP*, September 1993.
- [WTK86] A. Witkin, D. Terzopoulos, and M. Kaas. Signal matching through scale space. Proc. am. Assoc. Artif. Intel., pages 714-719, 1986.
- [WZ96] Hugh Williams and Justin Zobel. Indexing nucleotide databases for fast query evaluation. Proc. of 5-Th Intl. Conf. on Extending Database Technology (EDBT), pages 275-288, March 1996. Eds. P. Apers, M. Bouzeghoub, G. Gardarin.
- [Yao78] A. C. Yao. On random 2,3 trees. Acta Informatica, 9:159-170, 1978.
- [YGM94] Tak W. Yan and Hector Garcia-Molina. Index structures for selective dissemination of information under the boolean model. ACM TODS, 19(2):332-364, June 1994.
- [YMD85] N. Yankelovich, N. Meyrowitz, and N. Van Dam. Reading and writing the electronic book. *IEEE Computer*, pages 15-30, October 1985.
- [Zah71] C.T. Zahn. Graph-theoretical methods for detecting and describing gestalt clusters. *IEEE Trans. on Computers*, C-20(1):68-86, January 1971.
- [ZC77] A.L. Zobrist and F.R. Carlson. Detection of combined occurrences. CACM, 20(1):31-35, January 1977.
- [Zip49] G.K. Zipf. Human Behavior and Principle of Least Effort: an Introduction to Human Ecology. Addison Wesley, Cambridge, Massachusetts, 1949.
- [ZMSD92] Justin Zobel, Alistair Moffat, and Ron Sacks-Davis. An efficient indexing technique for full-text database systems. VLDB, pages 352-362, August 1992.

[ZRT91] P. Zezula, F. Rabitti, and P. Tiberio. Dynamic partitioning of signature files. ACM TOIS, 9(4):336-369, October 1991.

INDEX

Α

All pairs query, 58 Amplitude spectrum, 106 Amplitude, 99

В

Bit interleaving, 28 Black noises, 67 Blocks, 10 Bloom filters, 47 Boolean query, 19 Brown noise, 67

$\overline{\mathbf{C}}$

Cluster hypothesis, 47 Colored noises, 67 Column-orthonormal, 100 Combined indices, 20 Conjugate, 99 Cosine law, 87 Cosine similarity function, 48 Cosine similarity, 90 Cross-talk, 72

D

Data mining, 57 Database management systems, 7 Daubechies-4 Discrete Wavelet Transform, 115 DBMS, 7 Deferred splitting, 15, 35 Differential files, 47 Dimensionality curse, 22, 39, 73 Dimensionality reduction, 84 Discrete Fourier Transform, 66 Discrete Wavelet Transform, 113 Distance editing, 43 Division hashing, 12 DNA, 62 Dot product, 100

Е

Editing distance, 43, 59, 62, 83 Eigenvalues, 121 Eigenvectors, 121 Energy spectrum, 67, 106 Energy, 67 Exact match query, 19

F

False alarms, 46, 58 False dismissals, 58 False drops, 46 FastMap, 89 Fourier Transform Short Window, 113 Fractal dimension, 33, 35 Fractals, 116 Frequency leak, 108, 113

\mathbf{G}

Geographic Information Systems, 25

GIS, 25

н

Haar transform, 114 Hashing function, 12

SEARCHING MULTIMEDIA DATABASES BY CONTENT

Hashing, 11 dynamic, 13 extendible, 13 linear, 13 open addressing, 12 separate chaining, 12 spiral, 13 division, 12 multiplication, 12 Hermitian matrix, 100

Ι

154

Indices combined, 20 Inner product, 100

J

JPEG, 111

Κ

Karhunen-Loeve, 84, 123 Key-to-address transformation, 11 Keyword queries, 47

\mathbf{L}

Latent Semantic Indexing, 48, 130 Linear quadtrees, 27, 30 Lower-bounding lemma, 62 LSI, 48, 130

Μ

Main memory, 10 MBR, 34 Minimum bounding rectangle, 34 Multi-dimensional scaling, 85 Multiplication hashing, 12 Multiresolution methods, 116

Ν

Nearest neighbor query, 20, 26, 58 Noise black, 67 brown, 67 colored, 67 pink, 67

0

Oct-trees, 38 Orthonormal, 100

Ρ

PAMs, 21 Partial match query, 19 Phase, 99 Point Access Methods, 21 Point query, 26 Postings lists, 20 Power spectrum, 67, 106 Precision, 52

\mathbf{Q}

OBIC. 71 Quadtree blocks, 29 Quadtrees, 116 linear, 27, 30 Query by example, 71 Query Boolean, 19 exact match, 19 nearest neighbor, 20 partial match, 19 range, 19 sub-pattern match, 58 whole match, 58 all pairs, 58 nearest neighbor, 26, 58 point, 26 range, 26 spatial join, 26, 58 whole-match, 71

\mathbf{R}

Random walks, 67 Range query, 19, 26 Ranked output, 51 RDBMS, 7 Recall, 52 Regular expression, 42 Relation, 7 Relational model, 7 Relevance feedback, 51 Row-orthonormal, 101 WWW, 41

Z

Z-value, 28 Zipf's law, 44

\mathbf{S}

SAMs, 60 Secondary store, 10 Semi-joins, 47 Signature files, 45 Similarity function cosine, 48 document-to-cluster, 48 document-to-document, 48 Singular Value Decomposition, 48, 126Spatial Access Methods, 25, 60 Spatial join query, 26, 58 Spectrum, 106 SQL, 7Structured Query Language, 7 Superimposed coding, 46 SVD, 48

$\overline{\mathbf{T}}$

Text REtrieval Conference, 52 Time-warping, 83 Transpose matrix, 100 TREC, 52

v

Vector Space Model, 47

W

Wavelet transform, 113 Whole-match query, 71 World-wide-web, 41