

Matching and Retrieving Sequential Patterns Under Regression

Hansheng Lei, Venu Govindaraju
CUBS, Center for Unified Biometrics and Sensors
State University of New York at Buffalo
Amherst, NY 14260
{hlel, govind}@cse.buffalo.edu

Abstract

Sequential pattern matching and retrieving is of real value. For example, finding stocks in the NASDAQ market whose closing prices are always about \$ β_0 higher than or β_1 times as that of a given company. The problem reduces to linear pattern retrieval: given query X , find all sequence Y from database S so that $Y = \beta_0 + \beta_1 X$ with confidence C .

In this paper, we novelly introduce Simple Linear Regression (SLR) model to tackle this problem. We extend 1-dimensional R^2 to ER^2 for multi-dimensional sequence matching. In addition, we present the SLR+FFT pruning technique based on SLR to speed up retrieval without incurring any false dismissal. Experimental results on both synthetic and real dataset show that the pruning ratio of SLR+FFT can be above 99%. Applying the retrieval technique to real stocks discovered many interesting patterns. On-line signature recognition using ER^2 as similarity measure also demonstrated high accuracy.

1. Introduction

Sequence analysis has attracted a lot of research interests with a wide range of applications. A lot of work have been done in how to define a (dis)similarity measure which is fast, efficient and accurate (invariant to time warping, time scaling, amplitude scaling, offset shifting and etc.) [4, 13, 15, 16]. It is highly unlikely there exists any measure that can address all these challenges. Among the most popular and successful measures, Euclidean norm and Dynamic Time Warping (DTW) [2] are widely accepted and deeply explored. The former is faster than the latter, while the latter is more accurate than the former, because Euclidean distance is one-one matching but DTW allows many-one and one-many matching. However, neither Euclidean norm nor DTW is invariant to amplitude scaling and offset shifting. For example, we have two sequences Y , X

and $Y = \beta_0 + \beta_1 X$, where β_0 , β_1 are the offset shifting factor and the amplitude scaling factor respectively. If the we consider such Y and X have high similarity, the distance of either Euclidean norm or DTW will fail our expectation. To circumvent the problem of scaling and shifting, mean-deviation normalization is proposed, i.e., sequences are normalized by $(X - \bar{X})/\sigma(X)$, where \bar{X} is the mean and $\sigma(X)$ is the standard deviation of X . By this kind of normalization, shifting and scaling factors are gone. However, we can not just let these factors go. We should control and make use of them.

In some applications, we may need more information than just "similar". Here is an example: finding out stocks in the real market whose closing prices are always \$ β_0 higher than or β_1 times as that of a given company. The problem reduces to sequence retrieval: given query X , search sequence Y from database S so that $Y = \beta_0 + \beta_1 X$ with a certain confidence. Mean-deviation normalization is non-eligible to do this. It is tempting to think the scaling factor β_1 between Y and X is $\frac{\sigma(Y)}{\sigma(X)}$. However, it is not the case. Only under condition that Y and X have exact linear relation as $Y = \beta_0 + \beta_1 X$, we have $\beta_1 = \frac{\sigma(Y)}{\sigma(X)}$, otherwise, it does not hold. Exact linear relationship in random sequences such as stocks rarely exists. That is why we need some *approximation* or some kind of *confidence* during retrieval.

Linear Regression Model is commonly used in statistics and econometrics[11, 17]. Linear Regression has Simple Linear Regression, which models $Y = \beta_0 + \beta_1 X$, and Multiple Linear Regression, which models $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$. In this paper, we will introduce SLR for sequence matching and develop techniques for efficient retrieval of linear patterns from large databases.

The rest of the paper is organized as follows. We provide the background of SLR in section 2. Then, we extend 1-dimensional R^2 to ER^2 for multi-dimensional sequence matching. Next, we propose techniques for efficient retrieval under SLR. The proposed techniques are evaluated by experiments in section 5. Finally, section 6 draws conclusions.

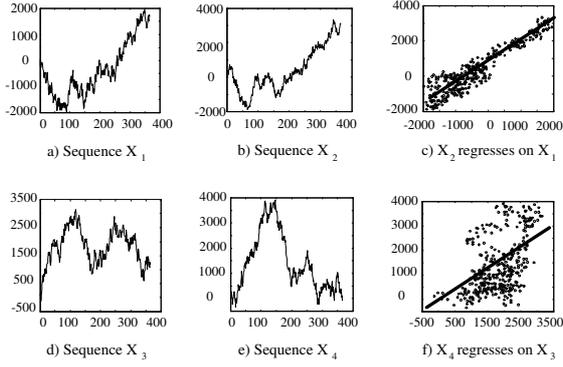


Figure 1. Two pairs of sequences. $R^2(X_1, X_2) = 0.91(91\%)$ and $R^2(X_3, X_4) = 0.31(31\%)$. X_1 and X_2 have high similarity (linearity) because the points in c) are distributed along a line. X_3 and X_4 have low similarity with scattered distribution in f).

2. Simple Linear Regression Background

Given two sequences $X = (x_1, x_2, \dots, x_n)$, $Y = (y_1, y_2, \dots, y_n)$, linear regression analysis is to statistically analyze the distribution of points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ in the X - Y space. If X and Y has strong linear relation, i.e., $Y \approx \beta_0 + \beta_1 X$, we can expect the distribution of these points is along a line, called the *regression line*. Fig. 1 shows examples.

To regress sequence Y on X , we first establish a *model*: $Y = \beta_0 + \beta_1 X + u$, where u is the *error term*. Note that $u = (u_1, u_2, \dots, u_n)$. Then we estimate the parameter β_0 and β_1 in the sense of minimum-sum-of squared-error, i.e., $\sum_{i=1}^n u_i^2 = \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2$ is minimized. From a geometric point of view, we estimate the regression line which is determined by β_0 and β_1 so that the line fits the points in the X - Y space as close as possible. Let $Q(\beta_0, \beta_1) = \sum_{i=1}^n u_i^2$. Note that Q is a function of β_0 and β_1 . To minimize $Q(\beta_0, \beta_1)$, we have $\frac{\partial Q}{\partial \beta_0} = 0$ and $\frac{\partial Q}{\partial \beta_1} = 0$. Starting from this, we can obtain the following results [11, 17]:

$$\beta_0 = \bar{Y} - \beta_1 \bar{X} \quad (1)$$

and

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2} \quad (2)$$

where $\bar{X} = \frac{1}{n} \sum_{i=1}^n (x_i)$, $\bar{Y} = \frac{1}{n} \sum_{i=1}^n (y_i)$.

With β_0 and β_1 as above, the regression line is determined. There remains a question: how well the regression line fits the points in the X - Y space? As a measure of the

goodness-of-fit, *R-squared* is defined as:

$$R^2 = 1 - \frac{\sum_{i=1}^n u_i^2}{\sum_{i=1}^n (y_i - \bar{Y})^2} \quad (3)$$

R-squared is also called the *coefficient of determination*. It can be interpreted as the *fraction of the variation in Y that is explained by X* . *R-squared* can be further derived as:

$$R^2 = \frac{[\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})]^2}{\sum_{i=1}^n (x_i - \bar{X})^2 \sum_{i=1}^n (y_i - \bar{Y})^2} \quad (4)$$

R^2 is exactly the square of Pearson's correlation r , i.e., the Cross-correlation:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{Y})^2}} \quad (5)$$

R^2 has following properties:

- **Reflexivity**, i.e., $R^2(X, X) = 1$.
- **Symmetry**, i.e., $R^2(X, Y) = R^2(Y, X)$. According to equation (4), no matter Y regresses on X or vice versa, R^2 is the same.
- $R^2 \in [0, 1]$. The closer the value to 1, the more the points tend to fall along the regression line. Thus, the stronger linear relation the two sequences have. $R^2 = 1$ means the two sequences have perfect linear relation, while $R^2 = 0$ means they have no linear relation at all.

Based on the properties above, R^2 is a good similarity measure. Fig. 1 show examples of two pairs of 1-dimensional sequences with length 400. It is clearly shown that high (low) R^2 value means high (low) similarity. Also, the similarity is a kind of confidence. Given two sequences, we can apply SLR and calculate β_0 , β_1 and R^2 . After that, we can say: *X and Y have linear relation as $Y = \beta_0 + \beta_1 X$ with confidence R^2 or X and Y have similarity R^2 .*

2.1. SLR vs. Mean-deviation Normalization

The R^2 has intrinsic relation with the normalized Euclidean distance by mean-deviation. Let $\nu(X)$ denotes the normalized version of X by mean-deviation, i.e., $\nu(X) = (X - \bar{X})/\sigma(X)$. Let $D_2(X, Y)$ denotes the Euclidean distance between sequence X and Y , then we have the following equation:

$$D_2(\nu(X), \nu(Y)) = \sqrt{2n(1 - r)}, \quad (6)$$

where r is the Pearson's r ; $R^2 = r^2$; n is the length of sequence. The proof of the equation is as follows.

$$\begin{aligned}
[D_2(\nu(X), \nu(Y))]^2 &= \sum_{i=1}^n \left(\frac{x_i - \bar{X}}{\sigma(X)} - \frac{y_i - \bar{Y}}{\sigma(Y)} \right)^2 \\
&= \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{\sigma(X)^2} + \frac{\sum_{i=1}^n (y_i - \bar{Y})^2}{\sigma(Y)^2} \\
&\quad - 2 \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sigma(X)\sigma(Y)} \\
&= n + n - 2nr = 2n(1 - r)
\end{aligned}$$

Note that $\sigma(X) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{X})^2$. Similar proof can be found in [18].

It is well known after the mean-deviation normalization, the shifting and scaling factors are discarded [9]. However, at the same time, the normalization loses the information of the shifting and scaling factors. From equations of (1) and (2), we can see that both β_0 and β_1 are more sophisticated than what we can obtain from the normalization. Only under *exact* linear relationship, i.e., $R^2 = 1$, we can derive $\beta_1 = \frac{\sigma(Y)}{\sigma(X)}$ by β_1^2/R^2 according to (2) and (4). However, exact linear relationship rarely exists in a domain with random distribution. This is why we need a distance tolerance or some kind of *confidence*. R^2 issued by SLR is a good confidence.

3. Extending R^2 to ER^2

Traditionally, SLR is only applied to 1-dimensional sequence. However, many kinds of sequences are multi-dimensional. For instance, the on-line handwritten signature sequence is multi-dimensional, including at least the x -, y -coordinates. To match a pair of K -dimensional sequences X and Y , we define ER^2 (Extended R-squared) as:

$$ER^2 = \frac{[\sum_{j=1}^K (\sum_{i=1}^n (x_{ji} - \bar{X}_j)(y_{ji} - \bar{Y}_j))]^2}{\sum_{j=1}^K \sum_{i=1}^n (x_{ji} - \bar{X}_j)^2 \sum_{j=1}^K \sum_{i=1}^n (y_{ji} - \bar{Y}_j)^2} \quad (7)$$

where \bar{X}_j (\bar{Y}_j) is the average of the j -th dimension of sequence X (Y).

It is trivial to prove that ER^2 has similar properties as R^2 , i.e., *reflexivity*, *symmetry* and $ER^2 \in [0, 1]$. ER^2 measures multidimensional sequences and directly tells how similar they are. In contrast, it is hard to normalize Euclidean distance or DTW distance to be intuitive.

The examples of similarity scores on on-line signatures in fig. 2 show the intuition behind ER^2 . Given only two sequences, R^2 or ER^2 directly tells how similar they are. Note that the distance based on Euclidean norm or DTW is relative [10]. If without comparison with other distances, the Euclidean or DTW distance itself does not mean any similarity. With the intuition, ER^2 can be utilized as the similarity metric in sequence classification problems. Fig.

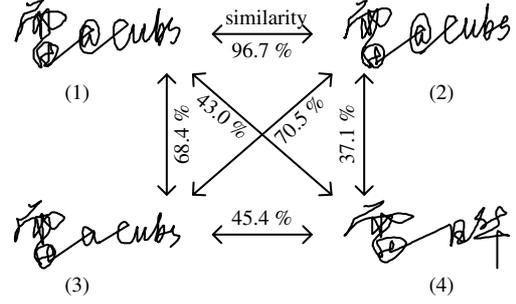


Figure 2. Four signatures and their similarities with each other defined by ER^2 . The similarity by ER^2 is much more intuitive than the dissimilarity distance, such as DTW or Euclidean distance.

ER^2						
	34.95%	98.42%	5.20%	36.28%	30.72%	44.35%
	0.44%	18.95%	96.75%	1.54%	00.06%	1.09%

Figure 3. Similarities by ER^2 between artificial 2D shapes.

3 shows examples on several artificial 2D shapes, where shapes are described by x -, y -coordinate sequences.

4. Retrieval under SLR

SLR can do more than "similar". Compared with regular range query like "given a query sequence, find sequences that are similar within distance ϵ " [1], the following query may interest users more.

Pattern Query *Select sequence Y , from database S , where $Y = \beta_0 + \beta_1 X$ with confidence $R^2 \geq C$ ($0 < C \leq 1$), $\beta_0 \in [\beta_{0l}, \beta_{0p}]$ and $\beta_1 \in [\beta_{1l}, \beta_{1p}]$. Pattern X , confidence threshold C , ranges of β_0 and β_1 are given by user.*

This query is more general. It reduces to regular query by fixing $\beta_0 = 0$ and $\beta_1 = 1$. It is also a kind of rule discovery problem. Rules such as "if stock X goes up, Y will go up β_1 times and Z will go down β_0 " can be expressed in linear relations.

Given a large database with massive objects, such as stocks, fingerprints, handwritten signatures, etc., it is very

expensive to search desired patterns with sequential scanning. That is why we need pruning skills to speed up retrieval. Can SLR support pruning? Fortunately, the answer is yes. The following theorems provide solid support.

Theorem 1. If sequence X and Y have linear relation as $Y = \beta_0 + \beta_1 X$ with confidence $R^2 \geq C$ ($0 < C \leq 1$) and $\beta_0 \in [\beta_{0l}, \beta_{0p}]$, $\beta_1 \in [\beta_{1l}, \beta_{1p}]$, then

$$\bar{Y} \in [\beta_{0l} + \min(\beta_{1l}\bar{X}, \beta_{1p}\bar{X}), \beta_{0p} + \max(\beta_{1l}\bar{X}, \beta_{1p}\bar{X})]$$

and

$$\sigma(Y) \in [\sigma(X)\min(|\beta_{1l}|, |\beta_{1p}|), \frac{\sigma(X)\max(|\beta_{1l}|, |\beta_{1p}|)}{\sqrt{C}}]$$

Proof. According to equation (1), $\bar{Y} = \beta_0 + \beta_1\bar{X}$. If $\beta_0 \in [\beta_{0l}, \beta_{0p}]$ (i.e., $\beta_{0l} \leq \beta_0 \leq \beta_{0p}$) and β_1 lies in the range $[\beta_{1l}, \beta_{1p}]$, \bar{Y} also lies in some range. It is not difficult to see this range is as stated in the above *theorem*.

$\sigma(Y)$ also must lie in some range. According to equation (2) and (4), we have $\frac{R^2}{\beta_1^2} = \frac{\sigma^2(X)}{\sigma^2(Y)}$. Thus, $\sigma^2(Y) = \frac{\beta_1^2 \sigma^2(X)}{R^2}$. With $R^2 \geq C$ and $\beta_1 \in [\beta_{1l}, \beta_{1p}]$, it is not difficult to conclude the proof. Note that β_1 can be negative, so, we take $\min(|\beta_{1l}|, |\beta_{1p}|)$ as the lower bound and $\max(|\beta_{1l}|, |\beta_{1p}|)$ as the upper bound.

Theorem 2. If sequence X and Y have linear relation as $Y = \beta_0 + \beta_1 X$ with confidence $R^2 \geq C$ ($0 < C \leq 1$) and $\beta_0 \in [\beta_{0l}, \beta_{0p}]$, $\beta_1 \in [\beta_{1l}, \beta_{1p}]$, then

1. If $\beta_{1l} \geq 0$, then

$$D_2(F_{2:k+1}(\nu(X)), F_{2:k+1}(\nu(Y))) \leq \sqrt{2n(1 - \sqrt{C})}$$

where $D_2(X, Y)$ denotes the Euclidean distance between X and Y , $\nu(X)$ is the normalized X by mean-deviation and $F_{2:k+1}(X)$ denotes the 2nd through the $(k+1)$ -th coefficients in the FFT (Fast Fourier Transform) domain of sequence X .

2. If $\beta_{1p} \leq 0$, then

$$D_2(F_{2:k+1}(\nu(X)), -F_{2:k+1}(\nu(Y))) \leq \sqrt{2n(1 - \sqrt{C})}$$

Proof. The difference in the two cases lies in that r (see equation (5)) is positive or negative.

1. $\beta_1 \geq \beta_{1l} \geq 0$ means $r \geq 0$, so, $r = \sqrt{R^2} \geq \sqrt{C}$. In addition, we know $D_2(\nu(X), \nu(Y)) = \sqrt{2n(1-r)}$. Therefore:

$$\begin{aligned} & D_2(F_{2:k+1}(\nu(X)), F_{2:k+1}(\nu(Y))) \\ & \leq D_2(F(\nu(X)), F(\nu(Y))) = D_2(\nu(X), \nu(Y)) \\ & = \sqrt{2n(1-r)} \leq \sqrt{2n(1-\sqrt{C})} \end{aligned}$$

where $F(X)$ denotes the FFT of X .

2. $\beta_1 \leq \beta_{1p} \leq 0$ means $r \leq 0$, so, $-r = \sqrt{R^2} \geq \sqrt{C}$. If $D_2(\nu(X), \nu(Y)) = \sqrt{2n(1-r)}$, then $D_2(\nu(X), -\nu(Y)) = \sqrt{2n(1-(-r))}$. The conclusion can be proved similarly as above.

The two theorems above support a speedup in retrieval. Suppose we have N raw sequences in the database S with length n . First step is to preprocess S and extract features. For each sequence Y in S , we create a tuple $[\text{ID}, \bar{Y}, \sigma(Y), F_{2:k+1}(\nu(Y))]$. ID is the identification number of sequence Y . Again \bar{Y} and $\sigma(Y)$ are the mean and standard deviation of Y respectively. Note that every coefficient in the FFT domain is a complex number, which has the real part and the imaginary part. Therefore, the tuple length is $3 + 2k$. Together we have N such tuples.

When a query is issued, i.e., user specifies the query sequence X , confidence threshold C ($0 < C \leq 1$), ranges $[\beta_{0l}, \beta_{0p}]$ and $[\beta_{1l}, \beta_{1p}]$, we do the following to retrieve qualified pattern Y from S :

1. **SLR-pruning.** According to theorem 1, we can predicate \bar{Y} and $\sigma(Y)$ lie in some range. Any candidate whose mean or standard deviation is not in the predicated range is pruned away immediately. If no candidate survives, then stop and return "no answer satisfying the query", otherwise, goto next step.
2. **FFT-pruning.** Extract FFT feature $F_{2:k+1}(\nu(X))$. Compare its Euclidean distance with those of the candidates which survived the SLR-pruning. If the distance is not in the desired range (see theorem 2), disqualify the candidate. If no candidate survives, then stop and return "no answer satisfying the query", otherwise, goto the next step.
3. **Post selection.** Collect all the survived candidates and retrieve the corresponding sequences from the databases according to their IDs. Test each one whether its regression on query pattern X satisfies $R^2 \geq C$ and β_0, β_1 lie in the range specified by the query.

The procedure above combines of the SLR-pruning and the FFT-pruning (SLR+FFT-pruning, hereafter, for brevity). Fig. 4 illustrates the architecture of the retrieval algorithm. The set of candidates which survive pruning is called the *superset* and the final set after post selection is called the *answer set* [1].

It has been noted that the SLR+FFT pruning is based on 1-dimensional SLR. If the sequence samples are multi-dimensional, we can apply the SLR+FFT pruning on the each dimension separately.

4.1. Discussion and Related works

Both *no false dismissals* and *high efficiency* are the ideal goals of a retrieval algorithm. False dismissal

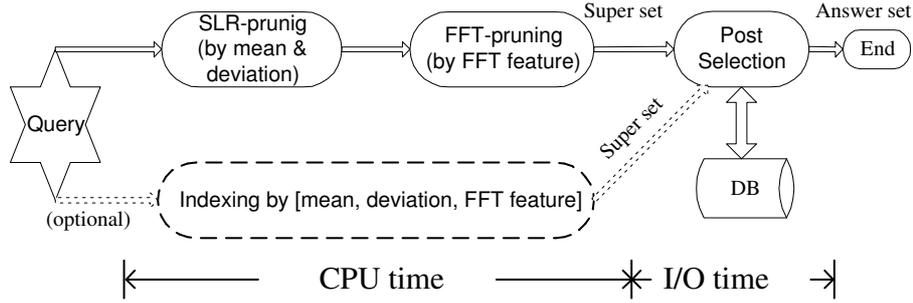


Figure 4. The flow chart of the retrieval algorithm.

means there exists qualified candidate missed during retrieval [1]. It may happen during pruning or approximating. The pruning techniques in our algorithm guarantee no false dismissals. First, the pruning by mean does not bring any false dismissal, because the qualified candidate must lie in range derived from user's query specification. Discarding all the candidates whose mean are out of this range speeds up the retrieval and guarantees the promising candidates remain in our scope. Second, the pruning by standard deviation also speeds up retrieval by disqualifying candidates directly. The pruning by FFT also introduces no false dismissals, which is guaranteed by Parseval's theorem [1]. The Euclidean distance between X and Y is same as between $FFT(X)$ and $FFT(Y)$. So, $D_2^2(FFT_{2:k+1}(X), FFT_{2:k+1}(Y)) \leq D_2^2(FFT(X), FFT(Y)) = D_2^2(X, Y)$. If even $D_2^2(FFT_{2:k+1}(X), FFT_{2:k+1}(Y))$ is beyond some threshold ε which we can not tolerate, the distance between raw X and Y must be non-tolerable. Therefore, we can discard Y without false dismissal (suppose the X is the query). By FFT transform, most of the energy of a sequence is concentrated in the first few coefficients. The first 2nd through the $(k + 1)$ -th coefficients are extracted as the feature vector. Note that the first FFT coefficient of the mean-deviation normalized sequence is always 0. Generally, $k \ll n$ and $k = 5$ is enough.

Pruning techniques is widely used in retrieval [8, 14, 3]. Actually, it is also implicitly used in indexing, which is to locate the candidates that qualify certain range(s) directly. The goal of pruning is to speed up retrieval by discarding disqualified candidates so that the I/O time to retrieve raw candidates in databases is reduced. It is well known that the I/O time dominates the total execution time. Therefore, before accessing to the physical databases, we need to determine a scope where the true answer(s) must lie in (*no false dismissal*) and the scope is small as possible (*high efficiency*).

R*-tree can be utilized to index the feature vector $FFT_{2:k+1}(\nu(Y))$ directly [1]. R*-tree has high performance in retrieving qualified vectorial features. Although we describe pruning skills in our algorithm, it is optional to

index $(\bar{Y}, \sigma(Y), F_{2:k+1}(\nu(Y)))$ directly by R*-tree or other spatial indexing structures. By SLR+FFT-pruning or indexing, the superset is the same. The only difference is that indexing may improve the CPU time, however, the dominating time is the I/O time.

Goldin et al. in [6] gave the ranges of mean and standard deviation under query $Q = aS + b$, where Q is query, S is the candidate and a, b have lower bounds and up bounds. However, the ranges are derived under the assumption that Q and S have *exact* linear relation as $Q = aS + b$. The ranges can not guarantee no false dismissals, since exact linear relation with query sequence rarely exists. Wong et al. in [9] determined the shifting and scaling factors for a geometrical point of view. The results of β_0 and β_1 are the same as SLR. However, they did not derive the intuitive measure R^2 , which can not be derived straightforwardly geometrically.

Indexing/retrieving sequential patterns under Euclidean norm [1, 5, 6, 9] or DTW [3, 14] has been extensively studied. They all depend on some (dis)similarity measure. However, to the best of our knowledge, little work has been done on the sequential pattern query problem based on SLR, which is intuitive and valuable in real applications.

5. Experiments

Our experiments focus on the following concerns:

- How about the power of SLR+FFT-pruning?
- Does the linear pattern retrieval make sense in the real applications such as mining stocks?
- How about the accuracy of ER^2 compared with DTW or Euclidean norm on handwritten sequences?

To answer the first question, we compare the *pruning ratios* of SLR-pruning, FFT-pruning and SLR+FFT-pruning on synthetic sequences as well as real on-line signature se-

Prun. ratio	SLR-pruning	FFT-pruning	SLR+FFT
2	85.11 %	97.84 %	99.42 %
4	84.81 %	97.57 %	99.30 %
6	84.56 %	97.15 %	99.02 %
8	83.48 %	96.84 %	98.76 %
10	82.48 %	96.03%	98.55%
Average	84.10 %	97.09%	99.01%

Figure 5. The pruning ratios of SLR, FFT and SLR+FFT when β_{1p} varies from 2 to 10 (and $\beta_{1l} = 1/\beta_{1p}$) with confidence threshold $C = 90\%$ and shifting factor β_0 in range $[-100, 100]$.

Prun. ratio	SLR-pruning	FFT-pruning	SLR+FFT
20	86.71 %	98.19 %	99.41 %
40	86.17 %	98.01 %	99.39 %
60	85.03 %	97.97 %	99.37 %
80	84.56 %	97.82 %	99.29 %
100	83.91 %	97.59%	99.20%
Average	85.28 %	97.92%	99.33%

Figure 6. The pruning ratios of SLR, FFT and SLR+FFT when β_{0p} varies from 20 to 100 (and $\beta_{0l} = -\beta_{0p}$) with confidence $C = 90\%$, scaling factor β_1 in range $[0.1, 10]$.

quences. The pruning ratio is defined as follows:

$$Pruning\ ratio = \frac{\# \text{ of pruned candidates}}{\# \text{ of all candidates in database}} \quad (8)$$

To answer the second question, we apply the algorithm to the real stocks to see whether we can mine out any interesting relations. For third concern, we compare the recognition accuracy of ER^2 with that of DTW and Euclidean distance on on-line signatures.

5.1. Experiments on Synthetic Sequences

Since stocks are intrinsically the same as random walk sequences [1] and due to lack of massive stocks databases, we use synthetic random walk sequences to evaluate the pruning power of SLR+FFT.

We generated 10 million sequences with length 1024. Each sequence $X = (x_1, x_2, \dots, x_{1024})$ is produced by random walk: $x_i = x_{i-1} + z_i$, where z_i is an independent, identically distributed (IID) random variable taken in range $[-100, 100]$. The first element x_1 is randomly taken from the range $[-100, 100]$.

The SLR-pruning has direct relation with the predicated ranges of the mean and standard deviation, which are further directly related to three parameters: confidence threshold C , range of β_0 and range of β_1 . Therefore, we recorded the pruning ratios by varying one of the three parameters and fixing the other two. We randomly generate a sequence and simulate a query. We repeat 100000 times and record the average of the pruning ratio for each combination of the three parameters.

The results of fig. 5 show that after SLR-pruning, averagely 84.10% of the 10 million candidates are disqualified immediately. Then, after the FFT-pruning, only 0.58%

Prun. ratio	SLR-pruning	FFT-pruning	SLR+FFT
60 %	81.80 %	76.57 %	94.24 %
70 %	83.09 %	84.43 %	96.31 %
80 %	84.48 %	90.93 %	97.29 %
90 %	85.16 %	97.97 %	99.35 %
100 %	88.78 %	100%	100%
Average	84.66 %	89.98%	97.44%

Figure 7. The pruning ratios of SLR, FFT and SLR+FFT when confidence C varies from 60% to 100% with shifting factor β_0 in range $[-100, 100]$ and scaling factor β_1 in $[0.1, 10]$.

(100%-99.42%) is remained! However, if we directly apply FFT-pruning without the SLR-pruning, 2.16% candidates survive, which means the I/O time will be about 3.7 (2.16/0.58) times as that of SLR+FFT. Since I/O time dominates the total execution time, saving on I/O time is a real improvement on the performance. From fig. 5 we can also see that the pruning power of SLR+FFT increases as the range of β_1 decreases because the ranges of both mean and standard deviation become narrower so that more candidates are pruned away. From fig. 6 we see the similar results.

Fig. 7 shows that when C increases from 60% to 100%, the pruning power increase rapidly. This can be explained like this: C is directly related to the internal Euclidean distance threshold ϵ (recall equation *theorem 2*). When user queries perfect linear relation ($C = 100$), the size of su-

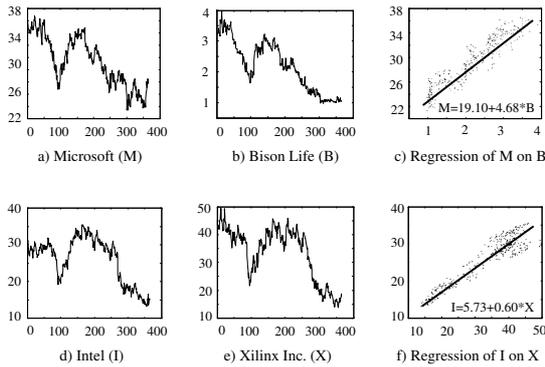


Figure 8. Two of the interesting linear relations retrieved out from the NASDAQ market.

perpet after SLR+FFT-pruning is almost zero, because exact linear relation rarely exists in random sequences.

In addition, we see that both SLR+FFT-pruning and FFT-pruning are relatively stable when confidence threshold C fixed, as show by fig. 5 and fig. 6. When C increases, both their pruning ratios increase rapidly, as shown by fig. 7. SLR-pruning is always stable, staying above 84%.

5.2. Experiments on Real Stocks

We downloaded all the NASDAQ stock prices from <http://table.finance.yahoo.com/> starting from May 8th, 2001 to May 8th, 2003. It has over 4000 stocks. We used the daily closing prices as sequences. These sequences are not of the same length for various reasons, such as cash dividend, stock split, etc. We made each sequence length 365 by truncating long sequences and removing short sequences. Finally, we have 3140 sequences all with the same length 365. The emphasis of the test is to observe whether we can discover interesting relations.

Since there are only 3140 sequences, we issued a query for each of them. The query was: given a query stock X , select stock Y from the dataset that satisfying $Y = \beta_0 + \beta_1 X$ with confidence $R^2 \geq 85\%$, $\beta_0 \in [-100, 100]$ and $\beta_1 \in [0.1, 10]$.

Fig. 8 shows two examples of the interesting relations we discovered in the stocks. The stock of Microsoft Corporation (M) has high linear relation with that of Boston Life Sciences Inc (B). The relation is $M = 19.10 + 4.68 * B$ with confidence $R^2 = 88.28\%$. The stock of Intel Corporation (I) and Xilinx Inc (X) also have strong relation as $I = 5.73 + 0.60 * X$ with confidence $R^2 = 88.81\%$. Such kinds of relations are intriguing and valuable for stock brokers.

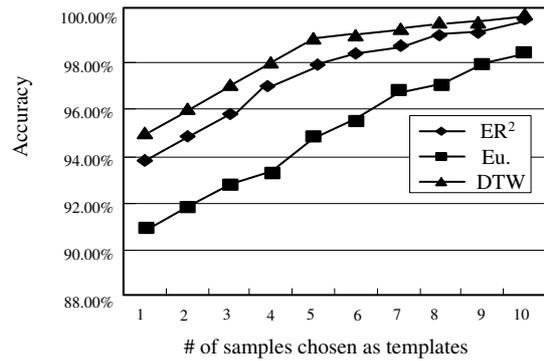


Figure 9. Accuracy of ER^2 , Euclidean norm and DTW on on-line signature classification.

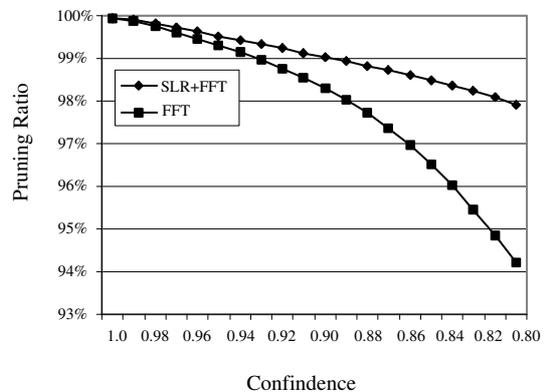


Figure 10. The pruning ratios of SLR+FFT and FFT on signature sequences when confidence varies C from 1.0 (100%) to 0.8 (80%).

5.3. Experiments on On-line Signature Sequences

The released SVC¹ on-line signature database was used to evaluate the accuracy of ER^2 as a similarity measure. Together, we have 80 classes of signatures with 20 genuine ones and 20 forgeries each class. So, totally, we had 3200 sequences. We preprocessed and uniformly *re-sampled* each signature to have the same length as 200 [12, 7]. Only the x -, y -coordinates were used to represent a signature, while other information such as pressure and time stamps were ignored. In this way, each signature is 200x2 sequence. We randomly chose M genuine signatures as templates from each class. Then, we randomly chose one signature to classify it by the 1-nearest

1 The First International Signature Verification Competition, 2004. <http://www.cs.ust.hk/svc2004/>.

neighbor. We repeated this procedure and calculated the average accuracy. We compared the accuracy by three (dis)similarity measures: ER^2 , DTW and Euclidean norm.

From the accuracies shown in fig. 9, We can see that ER^2 is a little bit less accurate than DTW, but more accurate than Euclidean norm. The higher accuracy of DTW comes from its elastic matching which is non-linear. However, DTW has much higher time complexity ($O(n^2)$) than ER^2 ($O(n)$). Considering the trade-off between accuracy and time, we believe that ER^2 has great potential for handwritten sequence matching. In our future work, we will try to combine ER^2 and DTW together to unify their advantages.

In addition, we tested the pruning ratio of SLR+FFT-pruning on the SVC signature database. For each signature, we issued a query to find all the similar signatures with similarity confidence R^2 no less than C . Note that the signature here was a 200x2 sequence. We actually issued two separate queries on x - and y -coordinates for each signature. For example, suppose $S_1 = [X_1, Y_1]$ is the query signature, then $S_2 = [X_2, Y_2]$ is retrieved out if $R^2(X_1, X_2) > C$ and $R^2(Y_1, Y_2) > C$. It also has to be noted that we implicitly fixed $\beta_0 = 0$ and $\beta_1 = 1$. For every possible C ranging from 80%-100%, we recorded the average pruning ratio. The results are shown in fig. 10. The pruning ratio of SLR+FFT is always higher than that of FFT alone. As threshold C decreases from 1 to 0.8 (80%), the SLR+FFT pruning outperforms FFT pruning more and more significantly. Although the signature database is not big enough, the high pruning ability of SLR+FFT is close to what we have seen on the synthetic sequences in section 5.1.

6. Conclusion and Future Work

We introduce SLR for sequential pattern matching and propose the SLR+FFT pruning technique for efficient sequence retrieval. In addition, similarity measure ER^2 is presented by extending 1-dimensional R^2 to multi-dimensional. FFT pruning is based on the Fast Fourier Transformation which concentrates the energy of a sequence to a few coefficients. SLR pruning is based on the linear regression which gives the predicted ranges of mean and standard deviation of a successful candidate. The combination of the two pruning techniques outperforms both SLR and FFT pruning alone. Experimental results show that the pruning ratio of SLR+FFT is high on both synthetic and real sequences. Simulation also shows that ER^2 is more accurate than Euclidean norm in multi-dimensional sequence classification, but less accurate than DTW. Our future work will study the possibility of ER^2 coupled with DTW for sequential pattern classification, such as on-line signatures, 2D shapes.

References

- [1] R. Agrawal, C. Faloutsos, and A. Swami, *Efficient similarity search in sequence database*, The 4th International Conference on Foundations of Data Organizations and Algorithms (1993), 69–84.
- [2] D. Berndt and J.Clifford, *Using dynamic time warping to find patterns in time series*, Working Notes of the Knowledge Discovery in Databases Workshop (1994), 359–370.
- [3] C. Faloutsos B.K. Yi, H.V. Jagadish, *Efficient retrieval of similar time sequences under time warping*, The 14th International Conference on Data Engineering (1998), 23–27.
- [4] Y. Manolopoulos C. Faloutsos. M. Ranganathan, *Fast subsequence matching in time-series databases*, (1994), 419–429.
- [5] A. Mendelzon D. Rafiei, *Efficient retrieval of similar time sequences using dft*, The 5th International Conference on Foundations of Data Organizations and Algorithms (FODO) (1998), 69–84.
- [6] D.Q. Goldin and P.C. Kanellakis, *On similarity queries for time-series data: Constraint specification and implementation*, The 1st International Conference on the Principles and Practice of Constraint Prog. (1995), 137–153.
- [7] A. Jain, F. Griess, and S. Connell, *On-line signature verification*, Pattern Recognition **35** (2002), no. 12, 2963–2972.
- [8] E. Keogh, *Exact indexing of dynamic time warping*, The 28th International Conference on Very Large Data Bases (2002), 406–417.
- [9] M.H. Wong K.W. Chu, *Fast time-series searching with scaling and shifting*, The 18th ACM Symp. on Principles of Database Systems (1999), 237–248.
- [10] H. Lei and V. Govindaraju, *GRM: Generalized regression model for clustering linear sequences*, The 4th SIAM International Conference on Data Mining (2004).
- [11] F. Mosteller and J. W. Tukey, *Data analysis and regression: A second course in statistics*, Addison-Wesley, 1977.
- [12] V. Nalwa, *Automatic on-line signature verification*, Proceedings of the IEEE **85** (1997), no. 2, 213–239.
- [13] H.S. Sawhney K. Shim R. Agrawal, K.I. Lin, *Fast similarity search in the presence of noise, scaling, and translation in time-series databases*, The 21st International Conference on Very Large Databases (1995), 490–501.
- [14] S. Park S. Kim, *An index-based approach for similarity search supporting time warping in large sequence databases*, The 17th International Conference on Data Engineering (2001), 607–614.
- [15] J. Yoon C. Hsu S. Park, W. Chu, *An index-based approach for similarity search supporting time warping in large sequence databases*, The 16th International Conference on Data Engineering (2000), 23–32.
- [16] Z.M. Ozsoyoglu T.Bozkaya, N. Yazdani, *Matching and indexing sequences of different lengths*, (1997), 128–135.
- [17] J. Wooldridge, *Introductory econometrics: a modern approach*, second ed., South-Western College Publishing, 1999.
- [18] Y. Zhu and D. Shaha, *Statstream: Statistical monitoring of thousands of data streams in real time*, The 28th International Conference on Very Large Data Bases (2002).