

data Structure project phase1

May 23, 2017

Ali A. Taheri
Ehsan E. Vakhshoori

Contents

1	time series and time series motifs	2
2	difference between match and trivial match in a time series	3
3	subsequence motif in a time series	3
4	datasets for time series	5
5	algorithm to detect time series motifs	6
6	programs output for a data	8
7	How much time it takes for algorithm to find motifs?	10

1 time series and time series motifs

time series

Time is the most important factor which ensures success in a business. Its difficult to keep up with the pace of time.

A Time Series is a sequence $T = (t_1, t_2, \dots, t_n)$ which is an ordered set of n real valued numbers.

The ordering is typically temporal however other kinds of data such as color distributions shapes and spectrographs also have a well defined ordering and can fruitfully be considered time series for the purpose of indexing and mining.

For example:

- Monthly unemployment rates for the previous five years.
- Daily production at a manufacturing plant for a month.
- Decade-by-decade population of a state of the previous century.

Components of a time series

Trend

The long-term tendency of a series to increase or fall (upward trend or downward trend).

Seasonality

The periodic fluctuation in the time series within a certain period. These fluctuations form a pattern that tends to repeat from one seasonal period to the next one.

Cycles

Long departures from the trend due to factors others than seasonality. Cycles usually occur along a large time interval, and the lengths of time between successive peaks or troughs of a cycle are not necessarily the same.

Irregular movement

The movement left after explaining the trend, seasonal and cyclical movements; random noise or error in a time series.

time series motifs

Time series motifs are pairs of individual time series, or subsequences of a longer time series, which are very similar to each other. Time series motifs are repeated segments in a long time series that, if exist, carry precise information about the underlying source of the time series. Motif discovery in time series data has received significant attention in the data mining community since its inception, principally because, motif discovery is meaningful and more likely to succeed when the data is large.

2 difference between match and trivial match in a time series

match

Given a positive real number R (called range) and a time series T containing a subsequence C beginning at position p and a subsequence M beginning at q , if $D(C, M) \leq R$, then M is called a matching subsequence of C .

The first three definitions are summarized in Figure 2, illustrating a time series of length 500, and two subsequences of length 128.

trivial match

Trivial Match: Given a time series T , containing a subsequence C beginning at position p and a matching subsequence M beginning at q , we say that M is a trivial match to C if either $p = q$ or there does not exist a subsequence M beginning at q such that $\text{dist}(C, M) > R$, and either $q < p$ or $p < q$.

3 subsequence motif in a time series

Time series motifs are pairs of individual time series, or subsequences of a longer time series, which are very similar to each other.

The k^{th} -Time Series motif is the k^{th} most similar pair in the database D . The

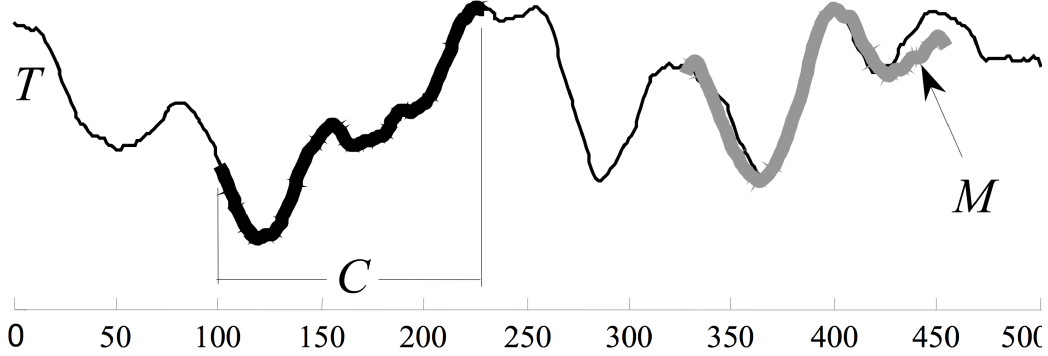


Figure 1: A visual intuition of a time series T (light line), a subsequence C (bold line) and a match M (bold gray line)

pair T_i, T_j is the k^{th} motif iff there exists a set S of pairs of time series of size exactly $k-1$ such that $\forall t_b$ and $\{t_j, t_d\} \in S$ and $\forall \{t_x, t_y\} \in S, \{t_w, t_b\} \notin S, dist(t_x, t_y) \leq dist(t_i, t_j) \leq dist(t_a, t_b)$.

Instead of dealing with pairs only we can also extend the notion of motifs to sets of time series that are very similar to each other.

The Subsequence Motif is a pair of subsequences $\{T_{i,n}, T_{j,n}\}$ of a long time series T that are most similar. More formally, $\forall a, b, i, j$ the pair $\{T_{i,n}, T_{j,n}\}$ is the subsequence motif iff $dist(T_{i,n}, T_{j,n}) \leq dist(T_{a,n}, T_{b,n}), |i - j| \geq w$ and $|a - b| \leq w$ for $w > 0$.

Note that we impose a constraint on the relative positions of the subsequence in the motif. This says that there should be a gap of at least w places between the subsequences. This restriction helps to prune out the trivial subsequence motifs. For example (and considering discrete data for simplicity), if we were looking for motifs of length four in the string: `sjdbbnvdfpgoeutyn-ABABABmbzchslfkeruyousjdq` Then in this case we probably don't want to consider the pair $\{ABAB, ABAB\}$ to be a motif, since they share 50% of their length (i.e AB is common to both). Instead, we would find the pair $\{sjdb, sjdq\}$ to be a more interesting approximately repeated pattern. In this example, we can enforce this by setting the parameter $w = 4$. In all the definitions given above we assumed there is a meaningful way to measure the distance between two time series. There are several such ways in the literature and our method is valid for any distance measure that is a metric. We use Euclidean distance and express it as $d(A, B)$. Recently extensive empirical comparisons have shown that the Euclidean distance is competitive with or superior to more complex measures on a wide variety of domains. Furthermore, its competitiveness increases as the datasets get larger, and it

is large datasets that are of interest here.

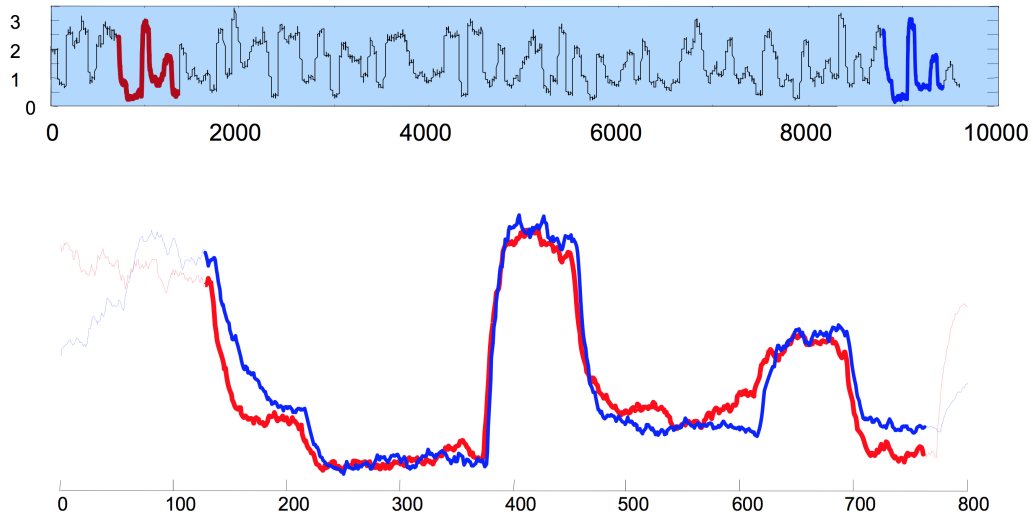


Figure 2: time motives

4 datasets for time series

Generate some input for main algorithm with c++:

(signal):

```
1 #include <iostream>
2 #include <ctime>
3 #include <cstdlib>
4 #include <fstream>
5 using namespace std;
6
7 int main()
8 {
9     freopen("input1.txt", "w", stdout);
10    srand(time(NULL));
11    for(int i=0; i<1000; i++)
12    {
13        for(int j=0; j<100; j++)
14            cout<<rand()%10<<" ";
15        cout<<endl;
16    }
17    return 0;
18 }
```

(string):

```
1 #include <iostream>
2 #include <ctime>
3 #include <cstdlib>
4 #include <fstream>
5 using namespace std;
6
7 int main()
8 {
9     freopen("input2.txt", "w", stdout);
10    srand(time(NULL));
11    for(int i=0; i<1000; i++)
12    {
13        for(int j=0; j<100; j++)
14            cout<<(char)(rand()%26+'a');
15        cout<<endl;
16    }
17    return 0;
18 }
```

5 algorithm to detect time series motifs

```
1 import time
2
3 def match(a, b, R):
4     count = 0
5     for i in range(len(a)):
6         if (type(a) is str):
7             if (a[i] != b[i]):
8                 count = count + 1
9         else:
10            count = count + abs(a[i] - b[i])
11    if (count <= R):
12        return True
13    else:
14        return False
15
16 def Algorithm_Find_Motif_Brute_Force(T, n, R):
17     target = open("output.txt", 'a')
18     for i in range(len(T) - n + 1):
19         count = 0
20         pointers = []
21         for j in range(i + n, len(T) - n + 1):
22             if (match(T[i:i+n], T[j:j+n], R)):
23                 count = count + 1
24                 pointers += [j]
```

```

25         target.write(str(T[i:i+n]))
26         target.write(" — ")
27         target.write(str(T[j:j+n]))
28         target.write("\n")
29     target.close()
30
31 target = open("output.txt", 'w')
32 target.close()
33 tm = open("outputTime.txt", 'w')
34 with open("input2.txt", "r") as ins:
35     arr = []
36     for line in ins:
37         arr.append(line)
38 for k in [2,5,6]:
39     start_time = time.time()
40     for i in range(1000):
41         string1=arr[i].split()
42         Algorithm.Find_Motif_Brute_Force(string1[0], k, 3)
43     tm.write(str(time.time()-start_time))
44     tm.write("\n")
45 tm.write("\n")
46
47 arr2=[]
48 with open("input1.txt", "r") as ins:
49     arr = []
50     for line in ins:
51         arr.append(line)
52 for i in range(len(arr)):
53     temp = arr[i].split()
54     arr2 += [[]]
55     for j in range(len(temp)):
56         arr2[i] += [ord(temp[j]) - 48]
57 for k in [2, 5, 6]:
58     start_time = time.time()
59     for i in range(len(arr2)):
60         Algorithm.Find_Motif_Brute_Force(arr2[i], k, 3)
61     tm.write(str(time.time()-start_time))
62     tm.write("\n")

```


6 programs output for a data

Input date: sjdbbnvdfpgoeutynABABABmbzchslfkeruyousjdq

output for this data with K=4 & R=1 :

```
1 sjdb — sjdq
```

output for this data with K=4 & R=2 :

```
1 sjdb — sjdq
2 euty — eruy
3 vnAB — ABAB
4 nABA — BABm
5 ABAB — ABmb
```

output for this data with K=4 & R=3 :

```
1 sjdb — vfdf
2 sjdb — ABmb
3 sjdb — slfk
4 sjdb — sjdq
5 jdbb — fdfp
6 jdbb — ABmb
7 jdbb — Bmbz
8 dbbn — dfpq
9 dbbn — tyvn
10 dbbn — Bmbz
11 dbbn — mbzc
12 bbnv — utyv
13 bbnv — yvnA
14 bbnv — mbzc
15 bbnv — bzch
16 bnvf — tyvn
17 bnvf — vnAB
18 bnvf — bzch
19 bnvf — hslf
20 nvfd — yvnA
21 nvfd — nABA
22 nvfd — slfk
23 nvfd — usjd
24 vfdf — vnAB
25 vfdf — hslf
26 vfdf — lfke
27 vfdf — sjdq
28 fdfp — slfk
29 fdfp — fker
30 dfpq — lfke
31 dfpq — sjdq
32 fpqo — fker
33 fpqo — ruyo
```

34 pqoe — lfke
35 pqoe — uyou
36 qoeu — fker
37 qoeu — keru
38 qoeu — uyou
39 qoeu — yous
40 oeut — keru
41 oeut — eruy
42 oeut — yous
43 oeut — ousj
44 euty — eruy
45 euty — ruyo
46 euty — ousj
47 utyv — ruyo
48 utyv — uyou
49 utyv — usjd
50 tyvn — uyou
51 yvnA — BABA
52 yvnA — yous
53 vnAB — ABAB
54 nABA — BABm
55 ABAB — ABmb
56 BABA — Bmbz
57 zchs — yous
58 chsl — ousj
59 hslf — usjd
60 slfk — sjdq

7 How much time it takes for algorithm to find motifs?

The brute force algorithm has a time complexity of $O(m^2)$.

This is critical because each distance computation takes $O(k)$ time which makes the brute force algorithm an $O(km^2)$ algorithm

