

Sweeping Incremental Algorithm for Matrix Profile (SIAMP)

Reza Akbarinia

Problem Definition

- Let
 - A and B : times series of size n
 - m: size of sequences
 - $A_m[i]$ or $B_m[i]$: sequence of size m starting at position i in A (or B)
 - $D_{i,j}$: Square of Euclidean distance between $A_m[i]$ and $B_m[j]$

Goal:

- Compute J_{AB} : such that $J_{AB}[i]$ returns the position of the nearest sequence of B to $A_m[i]$

Sweeping Incremental Algorithm for Matrix Profile (SIAMP)

- Main idea: compute the sequence distances incrementally
 - Each distance: in $O(1)$ instead of $O(m)$
 - Thus, an amortized complexity of $O(n)$ to find the nearest sequence of $A_m[i]$
- For this, we sweep the two time series in $n-m$ steps
 - In each step, the distances are incrementally computed, and minimum distances updated
 - In step k , we compute the distance of $A_m[i]$ and $B_m[i+k]$
 - i.e., sequences of A and B that have a difference of k in their initial positions

Algorithm

- For $i=0$ to $n-1$ $\text{Min_D}[i] := \infty$ //initialize minimum distances
- For $k=0$ to $n-m$ //sweep A and B in $n-m$ steps
 - Compute $D_{0,k}$ using Euclidean function
 - For $i=1$ to $n - k - 1$
 - Incrementally compute $D_{i, i+k}$ using $D_{i-1, i+k-1}$ // $O(1)$
 - If $(\text{Min_D}[i] > D_{i, i+k})$ then
 - $\text{Min_D}[i] := D_{i, i+k}$
 - $J_{AB}[i] := i+k$

Incremental Distance computation

- $D_{i,j}$: Square of Euclidean distance between $A_m[i]$ and $B_m[j]$
- $A_m[i] : \langle a_i, \dots, a_{i+m} \rangle$
- $B_m[j] : \langle b_j, \dots, b_{j+m} \rangle$
- $D_{i,j} = \sum (a_i - b_j)^2$ for $1 \leq i \leq m$
- $D_{i-1,j-1} = \sum (a_{i-1} - b_{j-1})^2$ for $1 \leq i \leq m$

Thus, we have

- $D_{i,j} = D_{i-1,j-1} - (a_{i-1} - b_{j-1})^2 + (a_{i+m} - b_{j+m})^2$

Analysis of SIAMP

- An exact algorithm for computing the matrix profile
- Time complexity: $O(n^2)$
- Space complexity: $O(n)$
- Simpler and faster than Keogh et al. algorithm whose complexity is $O(n^2 \log n)$
 - No need to Fourier transformations
 - No need to compute the mean and standard deviation of each sequence
 - No need to