Distance metrics for Sequential Data

- Dot matrix
- Sequence Alignment (Dynamic programming)
- Dynamic Time Wrapping (DTW)
- q-Gram distance

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (1)

What is sequential data ?

$$X = \left\{ x_1, x_2, \dots, x_N \right\} \ x_i \in \Omega$$

Sequence of measurement (s) of a quantity (or more) over time.

There is no need equally spaced time points

- > **Two types** of sequential data
 - Continuous time series (Ω real values)
 - Categorical (discrete) sequences (Ω alphabet)

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (2)

Continuous time series

- **Stock market prices** (for a single stock, or for multiple stocks).
- Heart rate of a patient over time.
- **Position** of one or multiple people/cars/airplanes over time (trajectories).
- **Speech:** represented as a sequence of audio measurements at discrete time steps.

• A **musical melody**: represented as a sequence of pairs (note, duration).

Categorical (discrete) sequences

- Biological sequences (DNA / Protein)
- Web navigation
- Song listening sequences
- Sunny/Rainy weather sequences
- Mobility in a network of cities / regions

Examples

Stock price (Bitcoin)

Daily degrees in a region





Human mobility



Dot matrix (Gibbs and McIntyre 1970)

Method for comparing two sequences to look for possible alignment

Algorithm for a dot matrix:

1. One sequence (A) is listed across the top of the matrix and the other (B) on the left side

2. Starting from the first character in B, one

moves across first row and placing a dot in many column where the character in A is the same

3. The process is continued in rows until all possible comparisons between A and B are made

4. Any region of similarity is revealed by a diagonal row of dots

5. Isolated dots not on diagonal represent random matches





Dot matrix

- Improve visualization of identical regions among sequences by using sliding windows, instead of writing down a dot for every character, that is common in both sequences
- We compare a number of positions (window size), and we write down a dot whenever there is minimum number of identical characters

two identical sequences

DNA1 on horizontal axis = 780 bases

DNA 2 on vertical axis = 780 bases



Click on plot to get positional data

two similar sequences sequences

DNA1 on horizontal axis = 1348 bases

DNA 2 on vertical axis = 1365 bases



Click on plot to get positional data

two very different sequences

DNA1 on horizontal axis = 1348 bases

DNA 2 on vertical axis = 2322 bases



Click on plot to get positional data

Sequence Alignment

Sequence alignment is a way of arranging the sequences to identify regions of similarity that may be a consequence of functional, structural or evolutionary relationships.



CGATGCAGACGTCA GATGCAAGACGTCA

The procedure of comparing two (pair-wise alignment) sequences is to search for a series of individual characters or patterns that are in the same order in the sequences.

Pairwise sequence alignment

Idea:

Display one sequence above another with **spaces (or gaps) inserted** in both to reveal similarity



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (12)

Alignment Scoring

- X= CTGTCG-CTGCACG
- Y= -TGC-CG-TG----
- Reward for matches: α Mismatch penalty: β Space penalty: γ

Score:
$$S(X,Y) = \alpha w - \beta x - \gamma y$$

w = #matches x = #mismatches y = #spaces

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (13)

Alignment Scoring

Reward for matches: Mismatch penalty: Space penalty:							•	10 2 5		
С	т	G	т	С	G	_	С	т	G	С
-	Т	G	С	-	С	G	-	Т	G	-
-5	10	10	-2	-5	-2	-5	-5	10	10	-5

Total similarity score = 11

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (14)

Types of Alignment algorithms

- Global alignment
- Local alignment
- Semi-Global alignment

CTGTCG-CTGCACG -TGC-CG-TG----

CTGTCGCTGCACG--

CTGTCGCTGCACG -TGCCG-TG----

- Requirements:
 - Scoring matrices
 - Gap penalty function

Optimum Alignment problem

- The score of an alignment is a measure of its quality
- Optimum alignment problem: Given a pair of sequences X and Y, find an alignment (global or local) with maximum score
- The similarity between X and Y, denoted sim(X, Y), is the maximum score of an alignment of X and Y.

• 2 sequences: X=(x₁, x₂, ..., x_n), Y=(y₁, y₂, ..., y_m)

- Suppose a similarity function S(a, b) between two members of discrete alphabet Ω
- Idea: Every edge in the directed graph has weight equal to the type of alignment.



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (17)

- 3 types of alignment
- Case 1: Characters from both sequences are aligned each other with cost S(x_i, y_i)



3 types of alignment

Case 2: Character from (column) sequence x is aligned with a gap from y with cost S(x_i, -) = - d
Y
Y
(i-1,j)
X
(i-1,j)
(i,j)

3 types of alignment

Case 3: Character from (row) sequence y is aligned with a gap from x with cost S(-, y_j) = - d
x
(i,i-1)
(i,j)

- Problem: Optimal path with the maximum total similarity cost
- Insert **Cost function** B(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial B(0,0) = 0



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (21)

- Problem: Optimal path with the maximum total similarity cost
- Insert **Cost function** B(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial B(0,0) = 0
- Case 1: alignment



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (22)

- Problem: Optimal path with the maximum total similarity cost
- Insert **Cost function** B(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial B(0,0) = 0
- Case 2



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (23)

- Problem: Optimal path with the maximum total similarity cost
- Insert **Cost function** B(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial B(0,0) = 0
- Case 3



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (24)

- Problem: Optimal path with the maximum total similarity cost
- Insert **Cost function** B(i,j), $0 \le i \le n$, $0 \le j \le m$

• Initial
$$B(0,0) = 0$$

• General step:

$$B(i, j) = \max \{B(i-1, j-1) + S(x_i, y_j), B(i-1, j) - d, B(i, j-1) - d\}$$

- Store which edge(s) was(were) selected
- At the end perform a back-tracking starting from B(n,m) to obtain the optimal path.

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (25)



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (26)



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (27)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (28)



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (29)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (30)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (31)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (32)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (33)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (34)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (35)



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (36)


Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (37)

Find the optimal path



Back-tracking

Each arrow introduces one character at the end of each aligned sequence.

• A <u>horizontal</u> move puts a gap in the <u>left</u> sequence.

• A <u>vertical</u> move puts a gap in the <u>top</u> sequence.

• A <u>diagonal</u> move uses one character from each sequence.



Back-tracking

Each arrow introduces one character at the end of each aligned sequence.

• A <u>horizontal</u> move puts a gap in the <u>left</u> sequence.

• A <u>vertical</u> move puts a gap in the <u>top</u> sequence.

• A <u>diagonal</u> move uses one character from each sequence.



2 possible alignments with the same cost (-6) Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (40)

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$

	Α	Α	Т	С	С	G	Α
Α							
С							
С							
Α							

Y=AATCCGA

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (42)

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$



Optimal path

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (43)

Y=AATCCGA

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$





Similarity score = - 2

 Find two common subsequences with maximum longest length and optimum cost of alignment

• If
$$\begin{aligned} X_k^i &= (x_k, x_{k+1}, \dots, x_i) & 1 \le k \le i \le n \\ Y_l^j &= (y_l, y_{l+1}, \dots, y_j) & 1 \le l \le j \le m \end{aligned} \text{ two subsequences} \\ \bullet \text{ Problem:} \qquad \begin{aligned} Find & \max_{\substack{1 \le k \le i \le n \\ 1 \le l \le j \le m}} \left\{ B(X_k^i, Y_l^j) \right\} \end{aligned}$$



• Totally there are $\binom{n}{2} \times \binom{m}{2}$ possible pairs of subsequences that must be examined

- **Dynamic programming** to solve this problem
- Introduce a score function L(i,j) and consider a (supposed) starting point I, where L(I)=0.
- L(i,j) denotes the maximum total score all paths starting from I and ending at (i,j).
- Then there are four (4) cases:

• There are four (4) cases:



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (48)

• Then there are four (4) cases:



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (49)

• Then there are four (4) cases:



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (50)

• Then there are four (4) cases:



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (51)

• There are four (4) cases:



Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (52)

- Problem: Longest Common Subsequence
- Cost function L(i,j), $0 \le i \le n$, $0 \le j \le m$

- Problem: Longest Common Subsequence
- Cost function L(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial L(0,0) = 0
- General step:

$$L(i, j) = \max \{L(i-1, j-1) + S(x_i, y_j), L(i-1, j) - d, L(i, j-1) - d, 0\}$$

- Problem: Longest Common Subsequence
- Cost function L(i,j), $0 \le i \le n$, $0 \le j \le m$
- Initial L(0,0) = 0
- General step:

$$L(i, j) = \max \{L(i-1, j-1) + S(x_i, y_j), L(i-1, j) - d, L(i, j-1) - d, 0\}$$

 At the end find the maximum L(i,j) and then perform a back-tracking until found zero (0). This will set the optimum local alignment.

Example 1 X=ACCA Y=AATCCGA

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$



Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (56)

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$

		Α	Α	Т	С	С	G	Α	
	0	0	0	0	0	0	0	0	$\left(L(i-1, j-1) + S(x_i, y_j)\right)$
Α	0	1	1	0	0	0	0	1	$L(i, j) = \max \begin{cases} L(i-1, j) - d \\ L(i, j-1) - d \end{cases}$
С	0	0	0	0	1	1	0	0	
С	0	0	0	0	1	2	0	0	
Α	0	1	1	0	0	0	0	1	

Scoring scheme:
$$S(a,b) = \begin{cases} 1 & a=b \\ -1 & a \neq b \end{cases}$$
, $d=2$

		Α	Α	Т	С	С	G	Α
	0	0	0	0	0	0	0	0
Α	0	1	1	0	0	0	0	1
С	0	0	0	0	1	1	0	0
С	0	0	0	0	1	2	0	0
Α	0	1	1	0	0	0	0	1

Local alignment (subsequence length 2)

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (58)



		Т	Т	С	Т	Α	Т	т	G
	0	0	0	0	0	0	0	0	0
Α	0	0	0	0	0	1	0	0	0
G	0	0	0	0	0	0	0	0	1
С	0	0	0	1	0	0	0	0	0
Т	0	1	1	0	2	0	1	1	0
Α	0	0	0	0	0	3	1	0	0
Α	0	0	0	0	0	1	2	0	0
С	0	0	0	1	0	1	0	1	0

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

$$L(i, j) = \max \begin{cases} L(i-1, j-1) + S(x_i, y_j) \\ L(i-1, j) - d \\ L(i, j-1) - d \\ 0 \end{cases}$$

		Т	Т	С	Т	Α	Т	Т	G
	0	0	0	0	0	0	0	0	0
Α	0	0	0	0	0	1	0	0	0
G	0	0	0	0	0	0	0	0	1
С	0	0	0	1	0	0	0	0	0
Т	0	1	1	0	2	0	1	1	0
A	0	0	0	0	0	3	1	0	0
Α	0	0	0	0	0	1	2	0	0
С	0	0	0	1	0	1	0	1	0

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

		Т	Т	С	Т	Α	Т	Т	G
	0	0	0	0	0	0	0	0	0
Α	0	0	0	0	0	1	0	0	0
G	0	0	0	0	0	0	0	0	1
С	0	0	0	1	0	0	0	0	0
Т	0	1	1	0	2	0	1	1	0
Α	0	0	0	0	0	3	1	0	0
Α	0	0	0	0	0	1	2	0	0
С	0	0	0	1	0	1	0	1	0

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

Local alignment (subsequence length 3)

X=CTA

Y=CTA

Semi Global Alignment

- 2 sequences: X=(x₁, x₂, ..., x_n), Y=(y₁, y₂, ..., y_m)
- In case where one sequence (e.g. x) has much larger length from the other (e.g. y), i.e n>>m.

Semi Global Alignment

- 2 sequences: X=(x₁, x₂, ..., x_n), Y=(y₁, y₂, ..., y_m)
- In case where one sequence (e.g. x) has much larger length from the other (e.g. y), i.e n>>m.
- Then, keep the **longer gapless** and allow gaps only to the smallest.

• Then: $B(i, j) = \max \begin{cases} B(i-1, j-1) + S(x_i, y_j) \\ B(i-1, j) - d \end{cases}$ (i-1,j-1) (i-1,j) (i-1,j

Semi Global Alignment

- 2 sequences: X=(x₁, x₂, ..., x_n), Y=(y₁, y₂, ..., y_m)
- In case where one sequence (e.g. x) has much larger length from the other (e.g. y), i.e n>>m.
- Then, keep the **longer gapless** and allow gaps only to the smallest.
- Then: $B(i, j) = \max \begin{cases} B(i-1, j-1) + S(x_i, y_j) \\ B(i-1, j) - d \end{cases}$ (i-1, j-1) (i-1, j) (i-
- At the end find the max B(i,j) value in the last column and traceback until 1st column.

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (65)

	Т	Α	С
Α			
G			
Α			
Т			
A			
Т			
С			
С			

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (66)

		Т	Α	С
	0	-2	-4	-6
Α	0	-3	-1	-5
G	0	-1		-2
Α	0	-1	0	-1
Т	0	1	0	-1
Α	0	0	2	-1
Т	0	1	0	1
С	0	-1	0	1
С	0	-1	-2	1

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (67)

		Т	Α	С
	0	-2	-4	-6
Α	0	-3	-1	-5
G	0	-1	-3	-2
Α	0	-1	0	-1
Т	0	1	0	-1
Α	0	0	2	-1
Т	0	1	0	1
С	0	-1	0	1
С	0	-1	-2	1

Scoring scheme:

$$S(a,b) = \begin{cases} 1 & a=b\\ -1 & a\neq b \end{cases}, \quad d=2$$

3 alignments (score=1)

X=AGATATCC Y=---TAC

X=AGATATCC Y=---TA-C-

Data Mining 2018 - Computer Science & Engineering, University of Ioannina - Sequential data (68)

DP Algorithm Variations







<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (69)

Dynamic Time Warping (DTW) distance

(Berndt & Clifford, 1994)

• Compare two time-series



- Use a distance function between two values, e.g. d(a,b)=(a-b)² or d(a,b)=|a-b|
- Following dynamic programming, we use a function
 D(i,j) that stores the minimum distance between the substrings (x₁, ..., x_i) and (y₁, ..., y_j). Then:

$$D(i, j) = d(x_i, y_j) + \min\{D(i-1, j-1), D(i, j-1), D(i-1, j)\}$$

- Not satisfying triangle inequality
- Slow to compute O(nm) time

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (70)

Euclidean vs. DTW distance

• Euclidean distance: Matches rigidly along the time axis. • DTW distance: Allows stretching and shrinking along the time axis.

<u>Data Mining 2018</u> – Computer Science & Engineering, University of Ioannina – Sequential data (71)

q - Gram distance

- 2 sequences: X=(x₁, x₂, ..., x_n), Y=(y₁, y₂, ..., y_m)
- q-gram: all the possible subsequences (w_q) of length q.
- There are $|\Omega|^q$ such possible subsequences.
- Then, q-gram distance:

$$D(X,Y) = \sum_{w_q} \left| n_X(w_q) - n_Y(w_q) \right|$$

where $n_X(w_q)$ frequency of w_q on sequence X

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (72)
q - Gram distance

Example: X = a a b a Y = a b a a $\Omega = \{a, b\}$

• q=2 : $w_a = \{aa, ab, ba, bb\}$. 2-gram distance

$$D(X,Y) = \sum_{w_q} \left| n_X(w_q) - n_Y(w_q) \right| = \left| 1 - 1 \right| + \left| 1 - 1 \right| + \left| 0 - 0 \right| = 0$$

q=3 : w_q = {aaa,aab,aba,abb,baa,bab,bba,bbb}.
3-gram distance

$$D(X,Y) = \sum_{w_q} \left| n_X(w_q) - n_Y(w_q) \right| = 0 + 1 + 0 + 0 + 1 + 0 + 0 = 2$$

Data Mining 2018 – Computer Science & Engineering, University of Ioannina – Sequential data (73)