# Dynamic Programming: Edit Distance

# Outline

- 1. DNA Sequence Comparison and CF
- 2. Change Problem
- 3. Manhattan Tourist Problem
- 4. Longest Paths in Graphs
- 5. Sequence Alignment
- 6. Edit Distance

# Section 5: Sequence Alignment

# Back to Biology: Sequence Alignment

- Recall that our original problem was to fit a similarity score on two DNA sequences:
- To do this we will use what is called an **alignment matrix**.

## Alignment Matrix: Example

• Given 2 DNA sequences **v** and **w** of length *m* and *n*:

v:ATCTGAT
$$m=7$$
w:TGCATA $n=6$ 

• Example of Alignment : 2 \* k matrix (k > m, n)



# Common Subsequence

• Given two sequences

 $v = v_1 v_2 ... v_m$  and  $w = w_1 w_2 ... w_n$ 

a common subsequence of  $\mathbf{v}$  and  $\mathbf{w}$  is a sequence of positions in  $\mathbf{v}$ :  $1 \le i_1 < i_2 < \ldots < i_t \le m$  and a sequence of positions in  $\mathbf{w}$ :  $1 \le j_1 < j_2 < \ldots < j_t \le n$  such that the  $i_t$ -th letter of  $\mathbf{v}$ is equal to the  $j_t$ -th letter of  $\mathbf{w}$ .

- Example:  $\mathbf{v} = \text{ATGCCAT}$ ,  $\mathbf{w} = \text{TCGGGCTATC}$ . Then take:
  - $i_1 = 2, i_2 = 3, i_3 = 6, i_4 = 7$
  - $j_1 = 1, j_2 = 3, j_3 = 8, j_4 = 9$
  - This gives us that the common subsequence is **TGAT**.

# Longest Common Subsequence

Given two sequences  $\mathbf{v} = v_1 v_2 \dots v_m$  and  $\mathbf{w} = w_1 w_2 \dots w_n$ 

the *Longest* Common Subsequence (LCS) of *v* and *w* is a sequence of positions in **v**:  $1 \le i_1 < i_2 < \ldots < i_T \le m$  and a sequence of positions in **w**:  $1 \le j_1 < j_2 < \dots < j_T \le n$ such that the *i<sub>t</sub>*-th letter of **v** is equal to  $j_t$ -th letter of **w** and T is maximal.

- Example:  $\mathbf{v} = \text{ATGCCAT}, \mathbf{w} = \text{TCGGGCTATC}.$ 
  - Before we found that TGAT is a subsequence.
  - The *longest* subsequence is **TGCAT**.
  - But...how do we *find* the LCS of two sequences?

• Assign one sequence to the rows, and one to the columns.



- Assign one sequence to the rows, and one to the columns.
- Every diagonal edge represents a match of elements.



- Assign one sequence to the rows, and one to the columns.
- Every diagonal edge represents a match of elements.
- Therefore, in a path <sup>A 4</sup> from source to sink, <sup>T 5</sup> the diagonal edges represent a common C<sup>7</sup> subsequence.





## Computing the LCS: Dynamic Programming

- Let  $\mathbf{v}_i$  = prefix of  $\mathbf{v}$  of length i:  $v_1 \dots v_i$
- and  $w_i$  = prefix of w of length j:  $w_1 \dots w_i$

The length of  $LCS(v_i, w_i)$  is computed by:

$$s_{i,j} = \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1 & if \ v_i = w_j \end{cases}$$

www.bioalgorithms.info

# Section 6: Edit Distance

# Hamming Distance

- The **Hamming Distance**  $d_H(\mathbf{v}, \mathbf{w})$  between two DNA sequences **v** and **w** of the same length is equal to the number of places in which the two sequences differ.
- Example: Given as follows,  $d_H(\mathbf{v}, \mathbf{w}) = 8$ :

V: ATATATATW: TATATATA

- However, note that these sequences are still very similar. ٠
  - Hamming Distance is therefore not an ideal similarity score, because it ignores insertions and deletions.

# **Edit Distance**

Levenshtein (1966) introduced the edit distance between two strings as the minimum number of elementary operations (insertions, deletions, and substitutions) needed to transform one string into the other

 $d(\mathbf{v}, \mathbf{w}) = MIN$  number of elementary operations to transform  $\mathbf{v} \rightarrow \mathbf{w}$ 

• So in our previous example, we can shift *w* one nucleotide to the right, and see that *w* is obtained from *v* by one insertion and one deletion:

V: ATATATATW: -TATATATA

- Hence the edit distance, d(v, w) = 2.
- Note: In order to provide this distance, we had to "fiddle" with the sequences. Hamming distance was easier to find.

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

## TGCATAT

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

#### **TGCATAT**(delete last **T**)

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

TGCATAT	(delete last T)
TGCATA	(delete last A)

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

TGCATAT	(delete last T)
TGCATA	(delete last A)
ATGCAT	(insert A at front)

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

TGCATAT(delete last T)TGCATA(delete last A)ATGCAT(insert A at front)ATCCAT(substitute C for G)

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

TGCATAT(delete last T)TGCATA(delete last A)ATGCAT(insert A at front)ATCCAT(substitute C for G)ATCCGAT(insert G before last A)

• We can transform TGCATAT  $\rightarrow$  ATCCGAT in 5 steps:

TGCATAT	(delete last T)
TGCATA	(delete last A)
ATGCAT	(insert A at front)
ATCCAT	(substitute C for G)
ATCCGAT	(insert G before last A)

• Note: This *only* allows us to conclude that the edit distance is *at most* 5.

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

# ATGCATAT

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

## ATGCATAT (insert A at front)

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

ATGCATAT	(insert A at front)
ATGCATAT	(delete second T)

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

ATGCATAT	(insert A at front)
ATGCA <mark>T</mark> AT	(delete second T)
ATGC <mark>G</mark> AT	(substitute G for A)

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

ATGCATAT(insert A at front)ATGCATAT(delete second T)ATGCGAT(substitute G for A)ATCCGAT(substitute C for G)

• Now we transform TGCATAT  $\rightarrow$  ATCCGAT in 4 steps:

ATGCATAT	(insert A at front)
ATGCATAT	(delete second T)
ATGCGAT	(substitute G for A)
ATCCGAT	(substitute C for G)

• Can we do even better? 3 steps? 2 steps? How can we know?

# Key Result

- **Theorem**: Given two sequences *v* and *w* of length *m* and *n*, the edit distance d(v,w) is given by d(v,w) = m + n s(v,w), where s(v,w) is the length of the longest common subsequence of *v* and *w*.
- This is great news, because it means that if solving the LCS problem for *v* and *w* is equivalent to finding the edit distance between them.
- Therefore, we will solve the LCS problem instead in the following slides.

## Return to the Edit Graph

• Every alignment corresponds to a path from source to sink.



## Return to the Edit Graph

- Every alignment corresponds to a path from source to sink.
- Horizontal and vertical edges correspond to indels (deletions and insertions).



## Return to the Edit Graph

- Every alignment corresponds to a path from source to sink.
- Horizontal and vertical edges correspond to indels (deletions and insertions).
- Diagonal edges correspond to matches and mismatches.



## Alignment as a Path in the Edit Graph: Example

- Suppose that our sequences are ATCGTAC, ATGTTAT.
- One possible alignment is:
- 0 1 2 2 3 4 5 6 A T \_ G T T A T \_ A T C G T \_ A \_ C 0 1 2 3 4 5 5 6 6 7
- Edit graph path:  $(0,0) \rightarrow (1,1) \rightarrow (2,2) \rightarrow$  $(2,3) \rightarrow (3,4) \rightarrow \text{etc.}$



#### Alignment with Dynamic Programming

Initialize  $O^{th}$  row and  $O^{th}$ column to be all zeroes.



#### Alignment with Dynamic Programming

- Initialize  $0^{th}$  row and  $0^{th}$  column to be all zeroes.
- Use the following recursive formula to calculate s<sub>i,j</sub> for each i, j:

$$S_{i, j} = max \begin{cases} S_{i-1, j-1} + 1 \\ S_{i-1, j} \\ S_{i, j-1} \end{cases}$$



- Note that the placement of the arrows shows where a given score originated from:
  - if from the top
  - $\leftarrow$  if from the left
  - $\bigvee$  if  $v_i = w_i$

• Continuing with the dynamic programming algorithm fills the table.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



- Continuing with the dynamic programming algorithm fills the table.
- We first look for matches, highlighted in red, and then we fill in the rest of that row and column.
- As we can see, we do not simply add 1 each time.



## Dynamic Alignment: Pseudocode

1. LCS(v,w)  
2. for 
$$i \in 1$$
 to  $n$   
3.  $s_{i,0} \in 0$   
4. for  $j \in 1$  to  $m$   
5.  $s_{0,j} \in 0$   
6. for  $i \in 1$  to  $n$   
7. for  $j \in 1$  to  $m$   
8.  $s_{i,j} \in \max \begin{cases} s_{i-1,j} \\ s_{i,j-1} \\ s_{i-1,j-1} + 1, \text{ if } v_i = w_j \\ s_{i-1,j-1} + 1, \text{ if } s_{i,j} = s_{i-1,j} \\ s_{i,j} \in 1 \\ s_{i,j} = s_{i,j-1} \\ s_{i,j} = s_{i-1,j-1} + 1 \\ s_{i,j} = s_{i,j} = s_{i-1,j-1} + 1 \\ s_{i,j} = s_{i-1,j-1}$ 

- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



- LCS(v,w) created the alignment grid.
- Follow the arrows backwards from the sink to the source to obtain the path corresponding to an optimal alignment:



# Printing LCS: Backtracking

```
PrintLCS(b,v,i,j)
1.
2. if i = 0 or j = 0
3.
            return
4. if b_{i,j} = " \smallsetminus "
             PrintLCS(b,v,i-1,j-1)
5.
6.
             print V<sub>i</sub>
7.
      else
            if b_{i,i} = " \uparrow "
8.
              PrintLCS(b,v,i-1,j)
9.
10.
            else
              PrintLCS(b,v,i,j-1)
11.
```

## LCS: Runtime

- It takes O(*nm*) time to fill in the *nxm* dynamic programming matrix: the pseudocode consists of a nested "for" loop inside of another "for" loop.
- This is a very positive result for a problem that appeared much more difficult initially.