

Multiple Sequence Alignment Using MATLAB

Meghna Mathur¹ and Geetika²

^{1,2}*Department of CSE/IT ITM University, Gurgaon, India.*

Abstract

Sequence alignment is an important task in bioinformatics which involves typical database search where data is in the form of DNA, RNA or protein sequence. For alignment various methods have been devised starting from pairwise alignment to multiple sequence alignment (MSA). To perform multiple sequence alignment various methods exist like progressive, iterative and concepts of dynamic programming in which we use Needleman Wunsch and Smith Waterman algorithms. This paper discusses various sequence alignment methods including their advantages and disadvantages. The alignment results of DNA sequence of chimpanzee and gorilla are shown.

Keywords: Multiple Sequence Alignment.

1. Introduction

In bioinformatics, a sequence alignment is a way of arranging the biological sequences including DNA (Deoxyribonucleic acid) or RNA (Ribonucleic acid) or protein. DNA sequencing is used to find genes, segments of DNA that code for specific protein or phenotype. Also used if a region has been sequenced it can be screened for characteristic features of genes.

Alignment can reveal homology between sequences. Similarity is the term that tells about the degree of match between two sequences. The principles include that sequence similarity do not always imply common function and conserved functions do not always imply similarity at sequence level.

Sequence alignments constitute an extremely powerful means of revealing the constraints imposed by the structure and function on the evolution of a protein and nucleic acid family, where as alignment task of multiple sequences requires large amount of computational time.

Local alignment is based on completeness. It includes the task of finding and extracting a pair of regions from two given biological sequences that exhibit high similarity. Global alignment is also based on completeness and is done across the entire sequence length to include as many matches as possible up to and including sequence end whereas Pairwise Sequence Alignment is used to identify regions of similarity that may indicate functional, structural and/or evolutionary relationships between two biological sequences (protein or nucleic acid). This type of alignment is based on numbers. Multiple sequence alignment (MSA) is the alignment of three or more biological sequences of similar length and therefore it is included in the alignment based on numbers. From the output of MSA applications, homology can be inferred and the evolutionary relationship between the sequences can be studied.

2. Methods of Alignment

2.1 Dot Matrix Method:

A dot matrix analysis is primarily a method for comparing two sequences to look for possible alignment of characters between the sequences. The method is also used for finding direct or inverted repeats in biological sequences and for predicting regions in RNA that is self-complementary and therefore have the potential of forming secondary structure through base-pairing. It works by locating regions of similarity between two sequences which provide a great deal of information about the function and structure of the query sequence. Similar structure indicates homology, or similar evolution, which provides critical information about the functions of these sequences. A dot matrix plot is a method of aligning two sequences to provide a picture of the homology between them. The dot matrix plot is created by designating one sequence to be the subject and placing it on the horizontal axis and designating the second sequence to be the query and placing it on the vertical axis of the matrix.

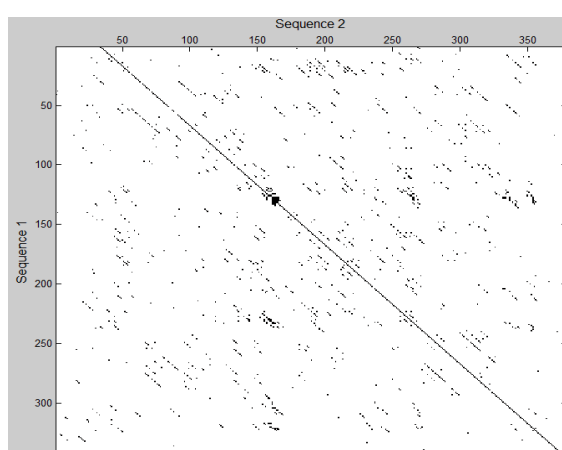


Figure 1: Sequence dot plot between Russian Neanderthal and German Neanderthal.

At each position within the matrix, a point is plotted if the horizontal and vertical elements are identical. Diagonal lines within the resulting matrix indicate regions of similarity. The advantages of dot matrix methods is that it readily reveals the presence of insertions/deletions and direct and inverted repeats that are more difficult to find by the other, more automated methods. However it has some shortcomings also. These are that most dot matrix computer programs do not show an actual alignment. Also it does not return a score to indicate how 'optimal' a given alignment is [10].

Example-Sequence1: Russian Neanderthal "AF254446"
Sequence2:German_Neanderthal "AF011222"

2.2 Dynamic Programming:

Dynamic programming (DP) algorithms are a class of algorithms typically applied to optimization problems. DP is applicable on an optimization problem with following key points. Firstly it should have an optimal substructure – an optimal solution to the problem contains within it optimal solutions to sub-problems. Secondly it must contain overlapping sub-problems i.e. the pieces of larger problem have a sequential dependency. DP works by first solving every sub-sub-problem just once, and saves its answer in a table, thereby avoiding the work of re-computing the answer every time the sub-sub-problem is encountered. Each intermediate answer is stored with a score, and DP finally chooses the sequence of solution that yields the highest score [10].

Both global and local types of alignments may be made by simple changes in the basic DP algorithm.

Scoring functions – example

w (match) = -2 or substitution matrix

w (mismatch) = -1 or substitution matrix

w (gap) = -3

Dynamic programming provides optimal alignment for a given set of scoring function which is its advantage. But it is slow due to the very large number of computational steps; computer memory requirement also increases as the square of the sequence lengths. Moreover the complexity is $O(n^2)$. Therefore; it is difficult to use the method for very long sequences. Dynamic programming has two algorithms that are used very frequently in sequence alignment Needleman Wunsch and Smith Waterman Algorithms.

2.2.1 Needleman Wunsch Algorithm for Global Alignment:

It performs a global alignment on two sequences. This algorithm is suitable when the two sequences are of similar length, with a significant degree of similarity throughout. It focuses on the best alignment over the entire length of the two sequences [1]. This algorithm works by first performing initialization where we create a matrix and fill its first row and column with the multiple of gap penalty then scoring takes place where in scoring of each row and column is obtained to get the final matrix and finally trace back is carried out in which we trace back to get the alignment of the respective sequences.

2.3 Progressive Methods:

Progressive methods are used for multiple sequence alignment. It is used to align three or more sequences. Using the standard dynamic programming algorithm on each pair, we can calculate the $(N*(N-1))/2$ (N is total number of sequences) distances between the sequence pairs. Distance matrix is obtained by applying the clustering algorithm and then constructing a guide tree. From the tree obtained we align the first node to the second node. After fixing the alignment, addition of another sequence or the third node takes place. After this we iterate the step until all the sequences are aligned. When a sequence is aligned to a group or when there is alignment in between the two groups of sequences, the alignment is performed that had the highest alignment score. The gap symbols in the alignment replaced with a neutral character [10].

2.4 Iterative Methods

These algorithms use iterative approach in which existing alignment can be realigned during addition of more sequences to multiple sequence alignment. Iterative methods can return to previously calculated pairwise alignments or sub-MSAs incorporating subsets of the query sequence as a means of optimizing a general objective function such as finding a high-quality alignment score [10]. Iterative methods are DIALIGN which takes an unusual approach of focusing narrowly on local alignments between sub-segments or sequence motifs without introducing a gap penalty. The alignment of individual motifs is then achieved with a matrix representation similar to a dot-matrix plot in a pairwise alignment. The most popular iteration-based method called MUSCLE (multiple sequence alignment by log-expectation) improves on progressive methods with a more accurate distance measure to assess the relatedness of two sequences. The distance measure is updated between iteration stages [6].

3. Sequence Alignment Tools

BLAST: Basic Local Alignment Search Tool, or BLAST, is an algorithm for comparing primary biological sequence information, such as the amino-acid sequences of different proteins or the nucleotides of DNA sequences. It addresses a fundamental problem and the heuristic algorithm it uses, is much faster than calculating an optimal alignment.

Using a heuristic method, BLAST finds similar sequences, not by comparing either sequence in its entirety, but rather by locating short matches between the two sequences.

FASTA is a DNA and protein sequence alignment software package which provides SSEARCH, an implementation of the optimal Smith-Waterman algorithm. Fasta format is a sequence in FASTA format begins with a single-line description, followed by lines of sequence data. The description line is distinguished from the sequence data by a greater-than (" $>$ ") symbol in the first column. Fasta algorithm FASTA ktups (k tuples) are shorter than BLAST words.

1-2 is for proteins and 4-6 is for nucleic acids. Lower ktups give a slower, more sensitive search. Higher ktups give a faster search with fewer false positives.

CLUSTALW Algorithm for Multiple Sequence Alignment where alignment can be done through different tools. CLUSTALW is one among the mostly accepted tool. CLUSTALW uses the progressive algorithm, by adding the sequence one by one until all the sequences are completely aligned. CLUSTALW follows following steps

- i. Determine all pairwise alignment between sequences and determine degrees of similarity between each pair.
- ii. Construct a similarity tree.
- iii. Combine the alignments starting from the most closely related groups to the most distantly related groups, using the “once a gap always a gap” rule.

The main disadvantage of this algorithm is that once a sequence has been aligned, that alignment can never be modified even if it conflicts with sequences added later in the process.

4. Evaluation of Sequence Alignment Algorithms

A sequence can be evaluated based on various factors like complexity of an algorithm, probability, accuracy and definiteness of an algorithm. Based on the complexity factor we can evaluate an algorithm. The algorithm that takes less time is more useful to the one that takes much time. This is so because our ultimate goal is always to solve a problem in as much less time as possible. In the case of ClustalW algorithm, the multiple sequence alignment of N sequences of length $O(n)$, the complexity is $O(N^2n^2)$. While by using the heuristics to this algorithm the complexity reduces to $O(N \log_2 N n^2)$. So it can be concluded that complexity helps in evaluating the use of an algorithm of a sequence alignment.

Second evaluating factor can be probability. Posterior probability is applied on pair of characters to find out the matching score. This helps in obtaining accurate results and higher speed [2].

The third evaluating factor can be accuracy of an existing algorithm. Accuracy can be defined as the correctness of an algorithm in terms of the output obtained on applying accurate inputs. An algorithm should always give one output to the number of inputs applied. So if that's the case then it can be said that the algorithm is accurate [3].

Fourth factor can be definiteness. By the term definiteness we mean that the algorithm should have finite number of steps. If an algorithm is not having a finite number of steps then that one cannot provide us with the correct results or the desired results.

5. Applications of Multiple Sequence Alignment

Multiple sequence alignment has emerged to have a lot of applications in the field of bioinformatics such as Sequence alignment helps in pattern recognition i.e. regions responsible for functional site can be identified by looking at the conserved regions,

profiles can be extracted using the file provided in multiple sequence alignment so that these can be used against databases, DNA regulatory elements can be located using multiple sequence alignment moreover tree reconstruction can take place by picking u related sequences which are used in multiple sequence alignment. Also decision about membership of protein to belong to a particular family can be inferred from multiple sequence alignment.

6. Result and Conclusion

Based on literature survey of all sequence alignment algorithms the table gives the summarized observation of all algorithms.

ALGORITHM	ADVANTAGES	DISADVANTAGES
Dot Matrix Method	It readily reveals the presence of insertions/deletions and direct and inverted repeats that are more difficult to find by the other, more automated methods.	Dot matrix computer programs do not show an actual alignment. It does not return a score to indicate how 'optimal' a given alignment is. Dot plots do not provide statistical analysis.
Dynamic Programming i) Needleman Wunsch Algorithm ii)Smith Waterman Algorithm	It is guaranteed in mathematical sense to provide an optimal alignment for a given set of scoring function. It does not require gap penalty.	It becomes slow as there are large computation steps. The memory requirement also increases as alignment sequences get large. It requires gap penalty for efficient working.
Progressive Alignment	These are efficient for aligning large set of sequences. Progressive alignment services are provided on web servers so the user need not install it locally.	Progressive alignment is not optimal globally. The errors generated at any stage are propagated to next stages which go till the end. When the sequences in the set are distantly related then the performance degrades.
Iterative Method i. DIALIGN ii. MUSCLE	They improve the accuracy of MSA. These can be repeated a number of times or until convergence.	

By studying various sequence alignment algorithms it can be concluded that progressive alignment is most widely used heuristic method for multiple sequence alignment. When posterior probability is applied to this method it helps in avoiding the errors occurring in early stages of alignment. The ant colony optimization method used in multiple sequence alignment helps in achieving higher speed of computation. ClustalW method used improves alignment by reducing time complexity and makes the aligning more dynamic.

References

- [1] Ankit Agrawal and Siddhartha Kumar Khaitan (2008), A New Heuristic for Multiple Sequence Alignment, *IEEE Iowa State University*, pp. 214-217.
- [2] Ling Chen, Welliu, Juan Chen (2007), Ant Colony Optimization Method for Multiple Sequence Alignment, *Intl. Conf. Machine Learning and Cybernetics*, Hong Kong, pp. 914-919.
- [3] Tristan Cazenve (2007), Overestimation for Multiple Sequence Alignment, in proceedings of the 2007 *IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology*, University Paris 8, pp.159-164.
- [4] C. Notredame (2002), recent progress in multiple sequence alignment: a survey, *Pharmacogenomics*, **vol. 3**, pp. 131-144.
- [5] Carsten Kemena, C.Notredame (2009), Upcoming Challenges for multiple sequence alignment methods in high-throughput area, *Oxford Journal*, **vol.25**, issue no.19, pp.2455-2465.
- [6] Robert C. Edger (2004), MUSCLE: multiple sequence alignment with high accuracy and high throughput, *Nucleic Acids Research*, **vol. 32**, pp. 1972-1977.
- [7] Mohamed Radhouene Aniba, Olivier Poch and Julie D. Thompson (2010), Issues in bioinformatics benchmarking: the case study of multiple sequence alignment, *Nucleic Acids Research*, **vol. 38**, pp.7353-7363.
- [8] EMBL-EBI, WellcomeTrustGenome Campus, Hinxton, Cambridge, CBIO, 1SD, UK [Online]
- [9] Available: <http://www.ebi.ac.uk/Tools/msa>.
- [10] David Mount, Bioinformatics sequence and genome analysis, cold spring harbor laboratory press, pp.238-260.
- [11] "Multiple Sequence Alignment :Progressive methods and HMMS", Terry Speed's computational Biology Course at UC Berkeley. [Online] Available: http://www-stat.stanford.edu/~nzhang/345_web/sequence_slides3.pdf