



Algorithms in Bioinformatics

CS 482



Bin Ma

Preface

Disclaimer Much of the information on this set of notes is transcribed directly/indirectly from the lectures of CS 482 during Winter 2022 as well as other related resources. I do not make any warranties about the completeness, reliability and accuracy of this set of notes. Use at your own risk.

For any questions, send me an email via <https://notes.sibeliusp.com/contact>.

You can find my notes for other courses on <https://notes.sibeliusp.com/>.

Sibeliusp Peng

Contents

Preface	1
1 Introduction	3
1.1 Objectives of this course	4
2 Sequence Alignment	5
2.1 Biology	5
2.2 Compare DNA sequences	6
2.3 Alignment	7
2.4 Score Function	8

Introduction

This course is officially titled as “Computational Techniques in Biological Sequence Analysis”. However, this is an old title and this course has been offered for over approximately twenty years. It should actually be called “Algorithms in Bioinformatics”. Around 1980, this area has a different name: Computational Biology. You may also hear another name: DNA computing. It is another area, and it does not solve biological problem. What is bioinformatics? It is biology + informatics. Biology is the reason, goal, purpose and informatics is the method.

Biology can be studied at different scales. In old days, biology tries to study organisms, namely living things, such as bacteria, animals. This is because people didn’t have tools to study at a lower level at that time. Now people study organs and tissues. Then people can look into cell level, molecular level. At molecular level, DNA is chain of nucleotide bases. Protein is chain of amino acids.

There are a lot of public and free molecular data. There are tremendous amount of public biomolecule data and free software. For example:

- NCBI’s sequence data bank: https://www.ncbi.nlm.nih.gov/nuccore/NC_045512
- PDB protein structure database: <https://www.rcsb.org/structure/6vxx>

Why do people do bioinformatics? The goal is to understand life at molecular level. Especially, for human health. For example, sequencing SARS-Cov2 genomes allowed people study the evolution of this virus. Study the structure of the spike protein and its interaction with the host cells. A lot of human diseases are related to genetics.

There are two areas of bioinformatics.

- Determine the molecule information, by analyzing the data produced by measuring instruments. Typically the data is in large scale and high throughput, which is hard for people to look at.
- Use the molecular data to make inference. For example, make prediction given the existing data.

Consider an example of genome sequencing. From [wikipedia](#),

The Human Genome Project (HGP) was an international scientific research project with the goal of determining the base pairs that make up human DNA, and of identifying, mapping and sequencing all of the genes of the human genome from both a physical and a functional standpoint. It remains the world’s largest collaborative biological project. Planning started after the idea was picked up in 1984 by the US government, the project formally launched in 1990, and was declared complete on April 14, 2003. Level “complete genome” was

achieved in May 2021.

Bioinformatics played an essential role in analyzing the data and assemble the genome. Today one can sequence a human's genome with $< \$1000$ in a couple of weeks. Bioinformatics is the key to utilize the NGS (next generation sequencing) data for genome sequencing. As such, today's cancer treatment starts to become personalized. And many new drugs now require gene sequencing as companion diagnostic.

1.1 Objectives of this course

- Know bioinformatics
 - Purpose and method
 - General topics
- Learn classic problems and algorithms in bioinformatics
- Learn wide-applicable computational techniques
 - String algorithms
 - Hidden Markov Model
 - Log likelihood ratio score
 - Statistical validation
 - A bit of machine learning

Some typical problems:

- Gene prediction problem
- Find the longest shared substring between human and mouse genomes. If we want to find similarities instead of exact matches, this will lead us to the homology search problem.
- Peptide Identification

Sequence Alignment

2.1 Biology

Consider two protein sequences:

- AVP78042.1 spike protein: MLFFL...
- YP_009724390.1 surface glycoprotein: MFVFL...

They look similar, how do we know these two proteins are similar? It's better to visualize them properly so that we can see the similarity. This is done by sequence alignment, and there are many existing tools: such as **Clustal Omega**. It's a common belief that these two proteins are developed from a common ancestor, then they evolve from an evolution tree. This is called homology.

There are two classes of nucleotide bases:

- Purine: A and G
- Pyrimidine: T and C

Base pairs are bonded by hydrogen bonds. Also, G-C bind stronger because of 3 H-bonds. DNA molecule is oriented.

So we know that DNA is double-helical, with two complementary strands. And the complementary bases: A-T, G-C. For example, the *reverse* complement of AAGGTAGC is GCTACCTT, because DNA is oriented.

Now consider DNA mutation. DNA mutates with a small probability when inherited by the offspring. For example, one base can be substituted by another, because there might be copy errors. This creates different alleles of the same gene. From wiki, allele is one of two, or more, forms of a given gene variant. When we inherit the DNA from parents, we only inherit half of each parent's genome. These together cause the differences between individuals of the same species.

Single Nucleotide Polymorphisms is a germline substitution of a single nucleotide at a specific position in the genome. Single base variation between members of a species. For Human, 90% of all human genetic variation is caused by SNPs. SNPs occur every 100 to 300 bases along the 3-billion-base human genome. It's a major risk for genetic disease, because when one base pair mutates, it will cause the express protein's functions.

2.2 Compare DNA sequences

The most often used distance on strings in computer science is Hamming distance. This makes some sense on comparing DNA sequences in some cases: substitution. But there are other mutations: insertion/deletion (indel), which cannot be modelled correctly by Hamming distance. Other DNA rearrangements can also happen. But substitutions and indel are the two mutations we concern the most for this course.

Edit distance

Instead, we can use **edit distance**. How “far” away are two sequences from each other? Edit distance is defined to be the minimum number of edit operations needed to convert one to another. Here edit operations include substitutions and indels.

Note that Edit distance is a distance metric:

- Identity: $d(x, y) = 0$ if and only if $x = y$.
- Symmetry: $d(x, y) = d(y, x)$.
- Triangular inequality: $d(x, z) \leq d(x, y) + d(y, z)$.

Now we prepare for the algorithm. For convenience of the proof, we treat each occurrence of the same letter different. For example, $ATAA \rightarrow ATA$, A can be done by either deleting the 2nd or 3rd letter A from the first string. These are different editing paths. This does not affect our definition of edit distance, but makes our later proof more precise.

Now it's ready to develop the dynamic programming algorithm for edit distance. Let $D[i, j]$ = edit distance between $S[1..i]$ to $T[1..j]$. Consider the edit operations associated with $S[i]$ and $T[j]$ the optimal edit operations. One of the following cases will happen:

1. $S[i]$ is deleted: $D[i, j] = D[i - 1, j] + 1$
2. $T[j]$ is inserted: $D[i, j] = D[i, j - 1] + 1$
3. $S[i]$ becomes $T[j]$: $D[i, j] = D[i - 1, j - 1] + \delta(S[i], T[j])$

where $\delta(S[i], T[j]) = 0$ if $S[i] = T[j]$ and 1 if not.

Algorithm 1: Dynamic Programming Algorithm for Edit distance

```

1  $D[0, 0] = 0$ 
2  $D[0, i] = i$  for  $i = 1..|S|$ 
3  $D[i, 0] = i$  for  $i = 1..|T|$ 
4 for  $i \leftarrow 1..|S|$  do
5   for  $j \leftarrow 1..|T|$  do
6      $D[i, j] = \min\{D[i - 1, j] + 1, D[i, j - 1] + 1, D[i - 1, j - 1] + \delta(S[i], T[j])\}$ 
7 return  $D[|S|, |T|]$ 
```

Longest Common Subsequence

Another way to evaluate the similarity of two sequences is through LCS. A subsequence is obtained by deleting some of the letters from the supersequence and concatenating the remaining letters together. For example, LCS of ATGCATTTA and ATGTACTTTC is ATGATTT. LCS can be computed with dynamic programming as well.

2.3 Alignment

The third way to compare to sequences is through sequence alignment. We want to insert spaces (-) to two sequences so that we can align them together and they are most similar column-wisely. For example,

```

ATGCA-TTTA
||| | || |
ATGTACTT-A

```

By “similar”, we usually need to use a scoring function. We define the alignment score to be **the total of column scores**. And each column is assigned by a constant score depending on matching conditions. For example, simple score scheme would be

- Match = 1
- Mismatch = -1
- indel = -1

Consider two alignments with the score scheme above:

```

AATGCGA-TTTT
|| | |||
G-TG--ACTTTC

```

has score 0.

```

AATG-CGATTTT
|| | ||
G-TGAC-TTTC-

```

has score -1.

As a side note, alignment can “simulate” LCS and edit distance. We have this following table for the scoring system:

	alignment	LCS	edit dist
match	+1	1	0
mismatch	-1	0	-1
indel	-1	0	-1

Now we can develop the dynamic algorithm for the alignment. Let $f(a, b)$ be the scoring scheme for a column with a and b . Here one of a and b can be the dash character -. Thus $f(-, x)$ and $f(x, -)$ represent scores of indels. The input is S, T . Let $D[i, j]$ be the optimal alignment score for $S[1..i], T[1..j]$.

Now we want to derive a recurrence relation. Suppose we are to align $S[1..i], T[1..j]$. Consider the *last column* of the optimal alignment. Three cases can happen.

$$\begin{array}{ccc}
 \begin{array}{cc} S[1..i-1] & S[i] \\ T[1..j-1] & T[j] \end{array} & \begin{array}{c} | \\ | \\ | \\ | \end{array} & \begin{array}{cc} S[1..i-1] & S[i] \\ T[1..j-1] & - \end{array} & \begin{array}{c} | \\ | \\ | \\ | \end{array} & \begin{array}{cc} S[1..i-1] & - \\ T[1..j-1] & T[j] \end{array} \\
 (1) & & (2) & & (3)
 \end{array}$$

Note that in each case, the sub-alignment without the last column is an optimal one. Prove by contradiction, assume there's a better sub-alignment without the last column. Then the optimal alignment (for the whole) should have used this better sub-alignment, instead of the current sub-alignment.

Then we have the recurrence for these three cases:

$$(1) D[i, j] = D[i - 1, j - 1] + f(S[i], T[j])$$

$$(2) D[i, j] = D[i - 1, j] + f(S[i], -)$$

$$(3) D[i, j] = D[i, j - 1] + f(-, T[j])$$

Algorithm 2: Sequence alignment

```

1  $D[0, 0] = 0$ 
2 for  $i \leftarrow 1..m$  do
3    $D[i, 0] = i \times \text{indel}$ 
4 for  $j \leftarrow 1..n$  do
5    $D[0, j] = j \times \text{indel}$ 
6 for  $i \leftarrow 1..m$  do
7   for  $j \leftarrow 1..n$  do
8      $D[i, j] = \max \begin{cases} D[i - 1, j - 1] + f(S[i], T[j]) \\ D[i - 1, j] + f(S[i], -) \\ D[i, j - 1] + f(-, T[j]) \end{cases}$ 
9 return  $D[m, n]$ 

```

Then the time complexity is $O(mn)$ where $|S| = m, |T| = n$. In practice, when we run the DP algorithm and build the DP table, we can add arrows from cell to cell: which of the three cases (sub-alignments) is chosen to get the current cell. Then we can backtrack the arrows and get the actual alignment. For backtracking, it is $O(n + m)$. Space complexity is $O(nm)$.

In practice, we don't need to record/store these arrows. Instead, when we backtrack, we can calculate the arrows based on the max of three options dynamically, which takes $O(m + n)$. This saves the cost of storing arrows, which costs $O(mn)$.

Moreover, if only score is needed, then space complexity can be reduced to linear space.

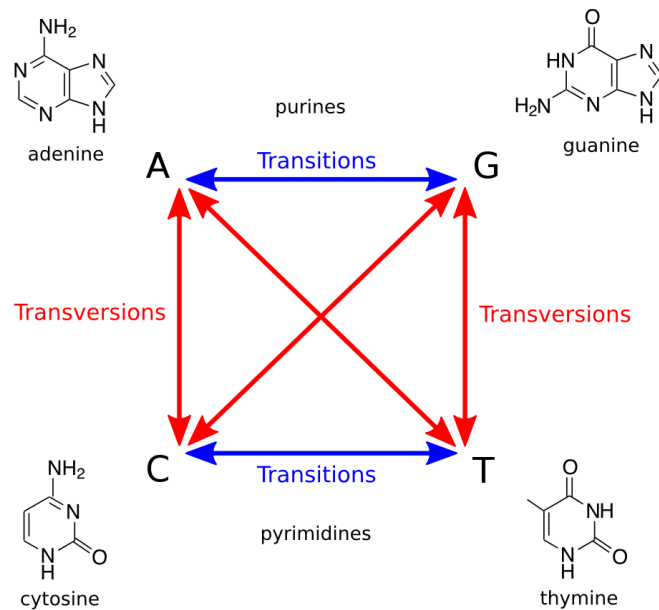
Note that the algorithm is designed for any score scheme $f(x, y)$. We indeed separate the algorithm and scoring. So we can optimize score scheme and the algorithm independently. Dijkstra once said:

The effective exploitation of his powers of abstraction must be regarded as one of the most vital activities of a competent programmer.

2.4 Score Function

Transition vs. Transversion

Recall that within DNA, there are two classes. One is purines (A, G), and the other is pyrimidines (C, T). Within the same class, they have similar chemical structure. Moreover, the mutation within the same class is called transitions; across the two classes, is called transversions. Transition is cheaper than transversions.



Picture from <https://dodona.ugent.be/en/activities/1351011626/>

Transition happens more frequently 2/3 of SNPs are transitions. In other words, transition is easier and therefore should be less penalized. For example,

AAAGCAAA	vs	AAAGCAAA
AAAT-AAA		AAA-TAAA

Observe that GT is transversion and CT is transition. The right alignment is probably better. This can be easily achieved by changing score scheme $f(a, b)$.

So how to build a score function? First, we need to know what you want. There are two purposes:

1. the optimal alignment reveals the true evolutionary history.
2. high score indicates homology (derived from same ancestor).

We want purpose 1 if possible, but purpose 2 is also useful.

For purpose 1, note that we might be wrong: score function might not have the power to reveal the history. For example, if the history goes as A → T → A, then in the sequence alignment, the final comparison is between A and A, which will not give us the true history. So you should never trust that an alignment must be right. It just optimizes the score. Should we give up purpose 1 at all?

For purpose 1, the optimal alignment may be *approximately correct under certain conditions* in practice. As long as we know the limitation, we can still use it. For example, for the following alignment, it is “very likely” the alignment is approximately equal to the evolutionary history.

```

ACGTATTACCGG-TTACCG
||| ||||| |||||
ACGGATTACCGGATTACCG

```

So we should keep in mind that when the score is low, alignment itself is not too useful.

Now let's try to improve our score function. Consider two alignment (with gaps):

AGATTTTTTTC	AGATTTTTTTTTTC
AGA---TTTTC	AGA-T-T-T-T-C

The left seems “simpler” than the right, intuitively. Indeed, during evolution, indels are relatively rare. However, insertion or deletion a segment of k consecutive bases is much easier than k scattered indels. But our current scoring method (adding up column scores) cannot distinguish the two. Currently, a gap of length k costs $k \times \text{indel}$. Thus, this is called the **linear gap penalty**. Denote the penalty function by $g(x)$, and it should have negative value. Left’s penalty is $g(3)$ and right is $3 \times g(1)$.

Consecutive insertions or deletions are called a gap. Suppose the gap penalty of a length k gap is $g(k)$ instead of the simple $c \times k$. Assume $g(x) + g(y) \leq g(x + y)$, thus a convex function. Otherwise does not serve the purpose of grouping indels. In this case, can the old DP algorithm still work? Recall the three cases:

$$\begin{array}{ccc}
 S[1..i-1] & S[i] & \\
 T[1..j-1] & T[j] & \\
 (1) & &
 \end{array}
 \quad
 \begin{array}{ccc}
 S[1..i-1] & S[i] & \\
 T[1..j-1] & - & \\
 (2) & &
 \end{array}
 \quad
 \begin{array}{ccc}
 S[1..i-1] & & - \\
 T[1..j-1] & T[j] & \\
 (3) & &
 \end{array}$$

First case is the same. The second case might be wrong. The last column might depend on previous columns. If the second to the last column of T is already a dash, then adding a new dash will increase the length of gap. The we cannot use DP algorithm anymore, in the sense that we cannot prove its correctness. We do not know the contribution of the last column to the gap penalty in the last two cases.

Assume we know the length of gap, then we can apply the recurrence relation. We still use $D[i, j]$ to denote the optimal alignment score of $S[1..i]$ and $T[1..j]$. We change cases 2 and 3 to include the last gap (not the last column). Then three cases change accordingly:

$$\begin{array}{ccc}
 S[1..i-1] & S[i] & \\
 T[1..j-1] & T[j] & \\
 (1) & &
 \end{array}
 \quad
 \begin{array}{ccc}
 S[1..i-k] & S[i-k+1..i] & \\
 T[1..j] & - \dots - & \\
 (2) & &
 \end{array}
 \quad
 \begin{array}{ccc}
 S[1..i] & & - \dots - \\
 T[1..j-k] & T[j-k+1..j] & \\
 (3) & &
 \end{array}$$

Then $D[i, j] = \max$ of the following three cases:

- (1) $D[i-1, j-1] + f(S[i], T[j])$
- (2) $\max_{1 \leq k \leq i} D[i-k, j] + g(k)$
- (3) $\max_{1 \leq k \leq j} D[i, j-k] + g(k)$

Now the time complexity will change: $O(mn(m+n))$, which is cubic.

In bioinformatics, very often we face the choice between:

- Reality: How close it approximates the real biology.
- Simplicity: How easy it can be computed.

We can simplify $g(k)$ a little bit. We basically want a function that grows slower than linear. We introduce **affine gap penalty**, in contrast to linear gap penalty:

$$g(k) = a + b \cdot k$$

where a is the gap open penalty: whenever we have a gap, we pay for it; b is the gap extension penalty: for example, once we open the gap, we pay less for the second indel within a group.

Example: Affine gap penalty

match = 1; mismatch = -1; gap open = -5; gap extension = -1.

```

ATAGG--AAG
|| || | |
ATTGGCAATG

```

6 match, 2 mismatch, 1 gap open, 2 gap extension, then the score is

$$6 - 2 + (-5 - 1 \times 2) = -3$$

```

ATAGG-AA-G
|| || || |
ATTGGCAATG

```

Similarly, the score is

$$7 - 1 - 5 - 1 - 5 - 1 = -6$$

Now the old algorithm doesn't work anymore. Consider the last column of an alignment again:

```

AT-GG-   ATGG--
|| ||    || |
ATTGGC   ATTGGC

```

When considering the last column, we need to know the second to the last column. When the last column is an indel, the added cost depends on the previous column. Because by induction we know the optimal solution for $D[i, j - 1]$, we can encode the previous column's configuration. We compute the optimal solution by limiting the last column to one of the following three configurations:

ATAGG	ATAGG-	ATAGGC
ATTGG	ATTGGC	ATTGG-
$D_0[i, j]$	$D_1[i, j]$	$D_2[i, j]$