Class 2 Lab*

Sequence Alignment & Database Searching (Pt. 1)

Barry Grant

Version 220926

i Instructions

Save this document to your computer and open it in a PDF viewer such as Preview (available on every mac) or Adobe Acrobat Reader (free for PC and Linux). Be sure to add your name and UC San Diego personal identification number (PID) and email below before answering all questions in the space provided.

Student Name UCSD PID UCSD Email
Richard Gao A16490010 r4gao@ucsd.edu

Overview:

Aligning novel sequences with previously characterized genes or proteins provides important insights into their common attributes and evolutionary origins.

In sections 1, 2 and 3 of this hands-on session we will first explore the principles and methods underlying the computational comparison and alignment of biomolecular sequences.

In section 4 we explore how these methods are used to search databases to identify homologues sequences (i.e. finding evolutionary related genes or proteins that are descended from a common ancestor).

In section 5 we highlight the detection limits of conventional BLAST. This sets the scene for introducing more sensitive (but often more time consuming) approaches including Profiles, PSI-BLAST and Hidden Markov Models (HMMs).

^{*}http://thegrantlab.org/teaching/

Section 1: Dot Plot Parameters

Dot plots are a simple graphical approach for the visual comparison of two sequences. They have a long history (see Maizel and Lenk 1981 and references therein) and entail placing one sequence on the vertical axis of a 2D grid (or matrix) and the other on the horizontal.

In its simplest form, a dot is placed where the horizontal and vertical sequence values match. More elaborate forms use 'sliding windows' composed of multiple characters and a threshold value, or 'match stringency' for two windows to be considered as matched.

Visit our very own simple dot plot web-app (Link1 or it's mirror Link2) and get a feel for how altering these major dot plot parameters change the displayed protein and DNA dot plots.

N.B. Note the questions listed on the web page (also found below) and add your answers in the space provided on the next page.

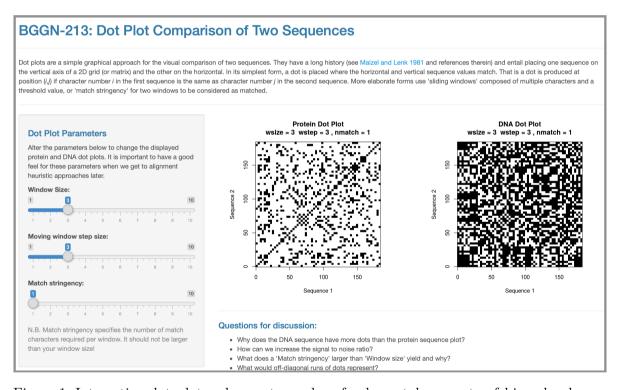


Figure 1: Interactive dot plot web app to explore fundemental concepts of biomolecular sequence comparasion

Q1 Why does the DNA sequence have more dots than the protein sequence plot?



Tip

What do you know about DNA composition vs protein composition?

DNA has 4 possible nucleotides while naturally occurring Amino Acids have 20 different possible identities so there is a greater chance that matches will occur between two DNA sequences versus two Amino Acid sequences.

Q2 How can we increase the signal to noise ratio?



? Tip

Signal in this case means correct matches that we actually want to highlight and noise means spurious matches that we don't want.

By first increasing the window size, which dictates the frame size of data points you are looking at, and subsequently increasing the match stringency which increases the number of match characters required per window. This means that data points will only show in which there are a sequence of matches which eliminates random data points where just 1 or 2 matches occur.

Q3 What does a 'Match stringency' larger than 'Window size' yield and why? Yields an error because match stringency is the number of match characters required per window, so the number of match characters required per window cannot be larger than the window size itself.

Q4 What are the major weaknesses of this approach?



? Tip

Is your inner nerd happy with this approach? How would you use it to determine if a second set of sequences was more similar to each other than a first set of sequences?

The dot plot compares data points 1:1 sequentially and doesn't account for potential gaps of matches in sequences or sequences of different lengths, so it cannot align sequences optimally. The dot plot is also a visual approach and doesn't provide a quantitative interpretation of similarity.

Section 2: Needleman-Wunsch Alignment

Sequence alignment methods often use something called a 'dynamic programming' algorithm that can be usefully considered as an extension of the dot plot approach. Here we have two sample sequences, and we'd like to use the Needleman-Wunsch algorithm discussed in class to align them. Feel free to use the clasroom white-boards and/or pen and paper and attach a photo to this PDF for gradescope.

Q5 Using a match score of +2, a mismatch score of -1, and a gap score of -2. Fill in the table below (or use pen and paper) for the following two sequences:

Sequence 1: ATTGC Sequence 2: AGTTC

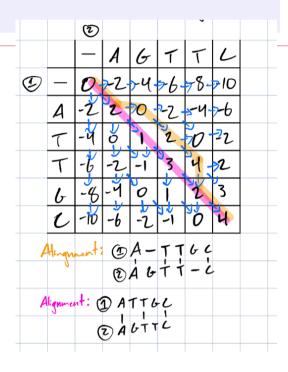
From your completed table what is the optimal score and corresponding alignment (with one sequence above the other)?



It can be hard to store the all important progress arrows in the PDF version of this document and thus you may prefer to use your own paper (or white-board) version that you can take a photo off for upload to gradescope.

The optimal alignment score is +4. Corresponding alignments in attached PDF.

		A	G	Т	Т	С
	0	-2	-4	-6	-8	-10
A	-2	2	0	-2	-4	-6
Т	-4	0	1	2	0	-2
Т	-6	-2	-1	3	4	2
G	-8	-4	0	1	2	3
С	-10	-6	-2	-1	0	4



Section 3: Practice makes perfect

Again use the Needleman-Wunsch algorithm discussed in class to align the following sequences:

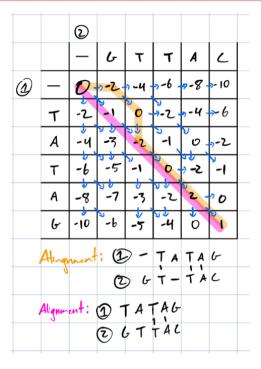
Q6 Using a match score of +2, a mismatch score of -1, and a gap score of -2. Fill in the table below and translate it into a alignment. What is the optimal score for this alignment? Is there one unique alignment with this score?

Sequence 1: TATAG

Sequence 2: GTTAC

The optimal score for this alignment is +1. There is more than one alignment with this score.

		G	Τ	Т	A	С
	0	-2	-4	-6	-8	-10
Т	-2	-1	0	-2	-4	-6
A	-4	-3	-2	-1	0	2
Т	-6	-5	-1	0	-2	-1
A	-8	-7	-3	-2	2	0
G	-10	-6	-5	-4	0	1



Section 4: Finding homologous sequence

Your collaborators found a protein while working on a fly species and have asked you to see if there are any human homologs.

>fly_protein

MDNHSSVPWASAASVTCLSLDAKCHSSSSSSSSKSAASSISAIPQEETQTMRHIAHTQRCLSRLTSLVAL LLIVLPMVFSPAHSCGPGRGLGRHRARNLYPLVLKQTIPNLSEYTNSASGPLEGVIRRDSPKFKDLVPNY NRDILFRDEEGTGADRLMSKRCKEKLNVLAYSVMNEWPGIRLLVTESWDEDYHHGQESLHYEGRAVTIAT SDRDQSKYGMLARLAVEAGFDWVSYVSRRHIYCSVKSDSSISSHVHGCFTPESTALLESGVRKPLGELSI GDRVLSMTANGQAVYSEVILFMDRNLEQMQNFVQLHTDGGAVLTVTPAHLVSVWQPESQKLTFVFADRIE EKNQVLVRDVETGELRPQRVVKVGSVRSKGVVAPLTREGTIVVNSVAASCYAVINSQSLAHWGLAPMRLL STLEAWLPAKEQLHSSPKVVSSAQQQNGIHWYANALYKVKDYVLPQSWRHD

Q7 Using the default settings for NCBI BLAST, can you find any homologs for this protein in Humans?



Tip

Try using the LIMITS and FILTERING options we covered in the last lab.

No

Q8 Try changing the database to refseq_protein. From the results, select a few proteins and find the common name for the species. What trend do you notice as you move down the results list?



Tip

Search google for the species name and use the taxonomy tab on your NCBI BLAST results page.

The common name for the species is Fly. The trend I notice is that as you go down the results, the e-value increases, the percent identity decreases, and other species begin to appear, specifically mosquitoes.

Q9 Finally, try also limiting the search to only H. Sapiens. What function do these proteins have?



You can simply type the Taxon ID 9606 in the "Organism" box.

These proteins participate in cell fate specification, cell-cell signaling, myelination, regulation of gene expression, and many other functions and processes in humans.

Q10 What function do you think this protein performs for your collaborators' organism? I think the protein would perform similar functions in our collaborators' flies as their functions in humans such as cell signaling, specification, and gene expression.

OPTIONAL EXTENSION

We will revisit this problem and introduce approaches with greater sensitivity (i.e. ability to find more remote homolouges) in the next lab.

Section 5: The limits of using BLAST for remote homologue detection

Let's return to the HBB protein that we explored in a previous class and see if we can find distantly related myoglobin and neuroglobin using HBB as a BLAST query.

>gi|4504349|ref|NP_000509.1| hemoglobin subunit beta [Homo sapiens] MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLG AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN ALAHKYH

After selecting **blastp** and entering the sequence, be sure to change the search database to "**refseq-protein**" and restrict our search organism to only **humans** (taxid: 9605). This will help focus our results to highlight distant homologs in humans.

Q11 What homologs did you find with this simple blastp search? Note their percent identities, coverage and E-values.

Hemoglobin subunit beta, 100% identity, 100% coverage, 3e-106 Hemoglobin subunit delta, 93.20% identity, 100% coverage, 1e-99 Hemoglobin subunit epsilon, 75.51% identity, 100% coverage, 2e-82

Now we could try changing the **Algorithm parameters** on the submission page to increase the number of hits reported. To do this you can click on the **Edit and Resubmit** link at the top left of your results page.

Q12 Try increasing the Expect threshold for your blasts search (e.g. to 2000). What new hits were reported? What about their alignment statistics? Do you trust these matches? Did you find myoglobin?

The new hits include different proteins such as "very long-chain specific acyl-CoA dehydrogenase" or "proto-oncogene tyrosine-protein kinase ROS isoform X14". Their percent identities and coverages are much lower, at around 40% and 20% respectively, and their E-values are much higher in the 2 to 3 hundreds. I do not trust these matches since the alignment statistics are so poor. I did not find myoglobin in the 92 sequences selected.

Q13: What one part of this exercise or associated lecture material is still confusing? If appropriate please also indicate the question number from this document and answer the question in the following anonymous form: Mudy_Point_Assessment_Form

Your comments will let us know which material needs to be further clarified and will help us gain stronger control of the material in this course. Thank you!

Discussion

Many useful 'rules of thumb' are expressed in terms of precent identity. If two proteins have more than 45% identical residues in their optimal alignment they typically have very similar structures and are likely to have a similar function. If two proteins have more than 25% identical residues (but less than 45% identity), they are likely to have a similar general folding pattern. Note that we will expand on the basis of this important sequence > structure > function relationship in the next lab.

Observations of a lower degree of sequence similarity cannot however rule out homology. Our very own late Russ Doolittle defined the region between 18-25% sequence identity as the "twilight zone" in which the suggestion of homology is tantalizing but dangerous. Below the **twilight zone** is a region where pairwise sequence alignments tell us very little - sometimes called the "midnight zone".

Our next class will introduce more advanced topics including profile and structure based approaches that can delve deeper into these important, but often hard to detect, sequence-structure-function relationships.