

Homology and Pairwise alignment exercises

Download the files that will be used

Browse to and download http://jbpease.github.io/crete2016/exercise_1.tar.gz

```
tar -xzf exercise_1.tar.gz
```

This should create a directory called exercise_1. You should be able to do everything in there

Working with pairwise alignments

```
python test.py
```

//you should see a NW pairwise alignment

Open test.py in a plain text editor. This is a python file that can run Needleman-Wunsch and Smith-Waterman in a verbose (but slow) way so you can experiment. The sc_mat is the scoring matrix. I have included EDNAFULL that is a simple DNA scoring matrix (used by NCBI) and I have included BLOSUM30, BLOSUM50, BLOSUM62, PAM50, PAM120, and PAM500 for amino acids (also used by NCBI).

Change the g (gap penalty) to -2

```
python test.py
```

Note not only how the alignment changes but how the matrix changes

Change the g (gap penalty) to -4

```
python test.py
```

Any changes?

Comment the current line (6) and uncomment the next line (7)

```
python test.py
```

How high does the gap penalty have to be to make it so there is only one gap?

Leave that line uncommented and uncomment the next line (8)

```
python test.py
```

Note the changes. The second one is a Smith-Waterman pairwise alignment of the same set of sequences.

Make the g = -10

```
python test.py
```

What changes?

If you change the gap penalty to 0 are the alignments the same? Similar?

Comment those two lines and uncomment the next two (9 and 10).

```
python test.py
```

Are these alignments similar or the same?

What penalty is required to make these different?

Now we will look at the protein alignments. Comment 9 and 10 and uncomment 12, 13 and 14.

```
python test.py
```

Browse to <ftp://ftp.ncbi.nih.gov/blast/matrices/> and see where you can download more PAM and BLOSUM. There are some already provided and you can change these to see if anything makes a difference.

Now we will look at `run.py`.

```
python run.py test1.fasta test2.fasta
```

We are going to examine what species these are later.

```
python run.py test1.fasta test3.fasta
```

How does this one look? Increase the gap penalty and try again.

What could be the reasons that this doesn't look as good. The script `revcomp.py` computes the reverse complement of a simple DNA sequence.

```
python revcomp.py test3.fasta test3.revcomp.fasta
```

```
python run.py test1.fasta test3.revcomp.fasta
```

Try this with a pairwise blast.

```
blastn -query test1.fasta -subject test3.fasta
```

```
blastn -query test1.fasta -subject test3.revcomp.fasta
```

How is blast letting you know that the first is reverse?

Working with SWIPE and BLAST

SWIPE is a Smith-Waterman tool for conducting pairwise Smith-Waterman alignments on a database of sequences. Some folks were unable to install the package so I will demonstrate it and then we will move to BLAST

The first thing we have to do for SWIPE or local BLAST is construct a special database based on a set of sequences that we want to compare. We can do that with the `makeblastdb` command. Everyone should have this.

```
makeblastdb -in rall.fasta -dbtype prot
```

Now we will run SWIPE on this database with the file

```
swipe -d rall.fasta -i r.fasta -v 1 -b 1
```

If we change the `b` to 0 it will just list the scores and the `e` values. This can be helpful for scanning results. Increase the `v` to get more results

```
swipe -d rall.fasta -i r.fasta -v 1 -b 0
```

```
swipe -d rall.fasta -i r.fasta -v 10 -b 0
```

What does it seem like this sequence is?

```
swipe -d rall.fasta -i r2.fasta -v 10 -b 0
```

What does it seem like this sequence is?

Now lets look at the actual alignments. Let's show 8.

```
swipe -d rall.fasta -i r.fasta -v10 -b 8
```

Is there anything interesting at this last sequences?

Now let's do a different file and see how that works.

```
swipe -d rall.fasta -i r3.fasta -v10 -b 8
```

Does this look like it is represented in the database?

We will now look at BLAST with the same set of files. Everyone should be able to do this part.

```
blastp -db rall.fasta -query r.fasta -num_alignments 0 -num_descriptions 1
```

This one might look like the cool one to begin with. To compare you (or I) can rerun the SWIPE with

```
swipe -d rall.fasta -i r.fasta -v 1 -b 1
```

Why might the scores be different?

Try this one with blast

```
blastp -db rall.fasta -query r3.fasta -num_alignments 0 -num_descriptions 1
```

When we did it with swipe, why was it different?

Let's work with some other files now.

```
makeblastdb -in r_p_3.fasta -dbtype prot
```

```
blastp -db r_p_3.fasta -query r3.fasta -num_alignments 0 -num_descriptions 1
```

What does it seem like this sequence is?

Look at r4.fasta. The others we have looked at are amino acid, what is this one? To use this, we need to use which blast?

```
blastx -db r_p_3.fasta -query r4.fasta -num_alignments 0 -num_descriptions 1
```

We can look at nucleotide blasts as well.

```
makeblastdb -in all.unall -dbtype nucl
```

```
blastn -db all.unall -query r4.fasta -num_alignments 0 -num_descriptions 1
```

We can also blast our protein sequence against the nucleotide database with tblastn

```
tblastn -db all.unall -query r3.fasta -num_alignments 0 -num_descriptions 1
```

Working with webBLAST

Often you will want to conduct blast analyses on the web before conducting or constructing pipelines that will use either BLAST or SWIPE

Go to <http://blast.ncbi.nlm.nih.gov/>

Open nerd.fasta in a plain text editor and copy it to the clipboard

Click on nucleotide blast and paste the sequence into the query sequence space. Then under the program selection, click on "somewhat similar", then click BLAST.

This is one of the more common exercises for BLAST, but it is fun so forgive me. This is the sequence used by Crichton in The Lost World. What do you notice about the sequence? What is it closely related to (top hit).

Go back and run a blastx on the same sequence. Scroll to the first few alignments. What do you notice in the gaps? If you are wondering Mark Boguski was a scientist at NCBI.

Conduct a blastn on the r4 sequence. Does the label match the sequence? Click on the accession for the first hit and you can see all the info for the sequence. Including some analyses.

If you go back you can see some interesting tools.

At the descriptions, click select all and then graphics. This gives a nice overview of the comparison of the sequences.

Go back and click on taxonomy reports at the top of the page. This will give you a rundown of the hits in taxa.

Go to the first blast page and go to the bottom and click on algorithm parameters. This will expand the options for the blast searches.

Change the word size to 15

match/mismatch to 1,-4

If you click open results in new window you can compare things more easily. Compare this with the one that you did before. (You can rerun with default parameters if you didn't leave them open).

Go to blastp and in the box type

>test

NIRVANA

Then click show results in new window and then click BLAST

Now go back and click algorithm, unclick "change parameters for short input sequences". Click BLAST again.

Now go back and click algorithm, make sure "change parameters for short input sequences" still off and increase "Expect threshold" to 100000 then click BLAST again. The "Expect threshold" is the number of hits you expect by random with larger being less stringent. Lower values are important for significance.

Conduct a blastp on r3.fasta but change the organism to Amborella. Is the first hit significant? This is another way to shrink the results but not the database. Compare the search summary of the first to a search summary without the organism limit. The statistics are the same so the size of the database didn't change.

Working with webBLAST if you have extra time

Go to blastp and try some different names that happen to overlap with the amino acid alphabet (ARNDCQEGHILKMFPSTWYVBZ). Note the evalues as you increase the letters.

You can try, NEIL, FLEA, ELVIS, SLASH, DARWIN, WALLACE, CRETE, ELLENKA

Report your best

Playing with significance

Nucleotide vs. Protein

We have a set of sequences that we can make a database of and then compare a nucleotide to that database and a protein to that database.

This is making a protein database (you can check the file in your text editor).

makeblastdb -in acorus_protein_gb.fasta -dbtype prot

Then we are blasting a file (r4p.fasta) against the newly made database. This is a protein file.

```
blastp -db acorus_protein_gb.fasta -query r4p.fasta -num_alignments 1  
-num_descriptions 1
```

Then we are blasting a file (r4.fasta) against the newly made database. This is a nucleotide file. Do you notice any differences? Which do you think might be more conservative (i.e., less likely to say they are similar)?

```
blastx -db acorus_protein_gb.fasta -query r4.fasta -num_alignments 1 -num_descriptions  
1
```

This is making a nucleotide database (you can check the file in your text editor) of the same sequences as above. We are then doing the blasts on this database (protein file against nucleotide database and nucleotide against nucleotide database).

```
makeblastdb -in acorus_nuc_gb.fasta -dbtype nucl
```

```
tblastn -db acorus_nuc_gb.fasta -query r4p.fasta -num_alignments 1 -num_descriptions 1
```

```
blastn -db acorus_nuc_gb.fasta -query r4.fasta -num_alignments 1 -num_descriptions 1
```

Database sizes

Here we have two files, a small and large. Check them out in your text editor. Then we will make databases of the two and blast our files against these databases to check whether they change the significance.

```
makeblastdb -in ITS.sm -dbtype nucl
```

```
makeblastdb -in ITS.all -dbtype nucl
```

```
blastn -db ITS.sm -query ITS.test -num_alignments 0 -num_descriptions 10
```

```
blastn -db ITS.all -query ITS.test -num_alignments 0 -num_descriptions 10
```

OK, I know that the results are all significant, but you should be able to see that the significance values change.