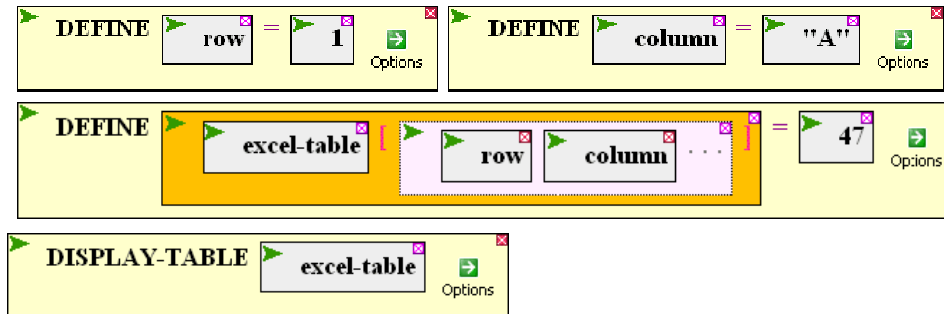


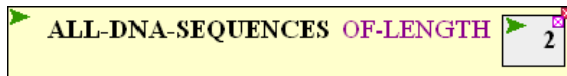
Introduction to Bioinformatics

Problem Set 5: Genome Analysis Investigations

1. Make and display an $x * y$ multiplication table for x and y going from 1 to 15. The [] function (in the List/Tables menu) will be essential here. You can use it to define elements of a two-dimensional table, as shown by example below:



2. Make and display a table containing information about the organisms known to BioBIKE. The information should include the name of the organism, the size of its genome, the number of genes, and its GC-fraction.
3. Determine the frequency of each dinucleotide in the genome of ss120. The following function will be useful:



4. Determine the number of times each codon is used in the coding genes of ss120.
5. Sort a list of 10 random numbers. Note that SORT lives in the LIST, List Production menu.
6. Make a list of lists, consisting of the following:
 - ((amino acid encoded by codon #1, codon #1, counts in ss120)
 - (amino acid encoded by codon #2, codon #2, counts in ss120)
 - ...)
7. Sort the above list and display it line by line. (Examine the options available for SORT)
8. What is the average size of the genes of *Prochlorococcus marinus* ss120? Of *Anabaena* PCC 7120? Are the genes of ss120 significantly smaller than the genes of A7120? Answer the question by doing a t-test. Note that there exists a T-TEST function, living in the MATHEMATICS / Statistics menu. Answer the question also by doing a simulation.
9. Return to *What is a gene*, part C, and construct by hand the table constructed in the investigation by MAKE-PSSM-FROM.
10. Construct a form that provides a palindrome twice the size of a sequence you give it. You'll want to know about the INVERSION-OF function in the STRING/SEQUENCE, String-production menu.

11. What is the frequency of 6-nt palindromic sequences in *Anabaena* PCC 7120?
- Construct a form that provides all palindromic sequences of length 6.
 - Count all palindromic sequences of length 6 in *Anabaena* PCC 7120.
 - Compare the incidence of each with the predicted number of counts.
 - Calculate for each the chi-squared score, comparing predicted with observed counts, then create and sort a list of palindromic sequences by its chi-squared score.
12. You have good reason to believe that the *Anabaena* gene *alr2328*, encoding glutamine synthetase is regulated by the availability of nitrogenous compounds. This makes sense, since glutamine synthetase incorporates ammonia into organic molecules and is a critical step in nitrogen metabolism. Every organism figures to have this gene, and it is probably regulated in the same way in every organism. With this in mind, you have some hope that you can find a regulatory sequence upstream from the gene by comparing the upstream sequences of all the orthologs of *alr2328* ("ortholog" means "gene related by evolutionary descent", more or less). See if you can find such a regulatory sequence, which would presumably serve as the binding site of a regulatory protein.
- You'll want to make yourself familiar with ORTHOLOGS-OF, in the GENES/PROTEINS menu, SEQUENCES-UPSTREAM-OF, in the GENES/PROTEINS, Gene neighborhood menu, and MOTIFS-IN, in the STRINGS/SEQUENCES, Bioinformatic-tools menu.
13. Construct a codon frequency table for the virus *myxococcus_phage_mx8* as described by Karlin (2001) [Trends in Microbiology 9:335-343]. In brief the table should have the following format (shown by example):

Codon-frequency["GGC"]

where the frequency for each codon is given as the *fraction* of times the codon is used from amongst all the codons encoding the same amino acid. Thus, the above example would be:

$$\frac{(\text{counts-of "GGC"})}{(\text{counts-of "GGA"} + \text{counts-of "GGC"} + \text{counts-of "GGG"} + \text{counts-of "GGT"})}$$

since GGA, GGC, GGG, and GGT all encode glycine.

Let's break up the task:

Create codon-count-table

- Generate a list of all codons in the genes of *myxococcus_phage_mx8*. You'll want to consider all the GENES-OF the phage, all the SEQUENCES-OF the genes, and then SPLIT the sequences every three nucleotides. Finally, you'll want to SIMPLIFY the resulting lists of lists ((...) (...) (...)...) into a single list.
- Define a **list** (call it `codon-count-list`) that contains the counts of each of the 64 triplet sequences amongst the codons of *myxococcus_phage_mx8*. Of course you'll want to make use of the COUNTS-OF and ALL-DNA-SEQUENCES functions, plus the result you obtained in **part a**.

- c. Learn how to create a single **element of a table** (call the table `codon-count-table`). Use `DEFINE` and `[]` functions to do this. For example:

`(DEFINE codon-count-table[any codon in quotes] AS any number of your choice)`

- d. Fill in `codon-count-table` with counts for all 64 triplet sequences. Use `APPLY-FUNCTION-OF`, using the function:

Function-of (codon codon-count)
= `(DEFINE codon-count-table[codon] = codon-count)`

and then apply that function of codon and codon-count to a list of codons and a list of codon counts. Note that you made a list of codon counts in part **a**, where you also used a list of codons. If your table has been defined properly, then you should be able to do things like this:

`codon-count-table["AAA"]`

and get the counts for that codon as the result (which should be 43).

Write a loop to calculate entries for `codon-frequency-table`

- e. Write a **loop** that considers in turn each of the 20 amino acids and displays the name of the amino acid
- f. Modify the loop so that you create a loop variable (call it **codons**) that you set equal (each iteration) to a list of codons that encode the specific amino acid. You may have to click the **reveal-all** option in the main menu of `FOR-EACH`. You'll be interested in a function `AA-TO-CODONS` that exists in the Genes-Proteins menu, Translation submenu. Display that list alongside the name of the amino acid.
- g. Modify the loop so that you create a loop variable (call it **codon-counts**) that you set equal (each iteration) to a list of codon counts for the codons you defined in **part f**. You'll want to make use of the function `ELEMENTS-OF-TABLE`. In place of *index* put the list of codons for the amino acid. Display the contents of this variable alongside the name of the amino acid and the list of codons.
- h. Modify the loop so that you create a loop variable (call it **count-sum**) that you set equal (each iteration) to the sum of the codon counts that you defined in **part g**. `SUM-OF` will be a useful function. Display the contents of this variable alongside all the other things you're already displaying.
- i. Modify the loop so that you create a loop variable (call it **codon-frequencies**) that you set equal (each iteration) to the codon-counts divided by the count sum. This should give you a list of frequencies that add up to one. Display this new variable alongside the others and confirm that the frequencies do add up to 1.
- j. Now that you have the frequencies, all that's left is to enter each of the frequencies into a table called **codon-freq-table** within the body of the loop. You'll do this in very much the same way as you created **codon-count-table** in **part d**.
- k. See what you got, using the `DISPLAY-TABLE`