

Homework 2

Simple Sevenless

Write a program `sevenless.py` to print out all the numbers from 0 to 99, one on each line, except, do not print any number perfectly divisible by 7.

Count up

Compute let's calculate some statistics from this GFF file which lists the location of genes and exons locations. Remember **GFF** is a structured format, tab delimited, which describes locations of features in a genome.

Here is a GFF file for the E. coli K-12 genome. ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/gff3/bacteria_0_collection/escherichia_coli_str_k_12_substr_mg1655/Escherichia_coli_str_k_12_substr_mg1655.ASM584v2.37.gff3.gz

Here is a FASTA file for the genome of E. coli K-12. ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_0_collection/escherichia_coli_str_k_12_substr_mg1655/dna/Escherichia_coli_str_k_12_substr_mg1655.ASM584v2.dna.chromosome.Chromosome.fa.gz

Write a script called `count_up.py` to:

1. Download this file (this can be in UNIX before you run your python script or you can incorporate this into the python). I already wrote part of this for you in the template code you can start with that executes a `curl` command from within your script. But if this doesn't make sense to you, you can remove that.
2. Count up and print out the number genes (gene feature)
3. Compute the total length of the genes (length is the `END - START`)
4. Use the FASTA file to compute the total length of genome (by adding up the length of each sequence in the file). Recall I lectured on a basic code to read in a FASTA file - you can also see that code template [here](#)
5. Print out the percentage of the genome which is coding

Codon compute

Use the following files to examine codon usage across these two bacteria. Remember that codons are triplets (eg ACA, GAT, ...). There are 64 total possible triplets. To count these, know that they are non-overlapping sets of three adjacent bases in the sequences, start with the very first base as the [reading frame](#).

These files are coding sequences of the predicted genes in each of two species.

1. ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_0_collection/salmonella_enterica_subsp_enterica_serovar_typhimurium_str_lt2/cds/Salmonella_enterica_subsp_enterica_serovar_typhimurium_str_lt2.ASM694v2.cds.all.fa.gz

2. ftp://ftp.ensemblgenomes.org/pub/bacteria/release-45/fasta/bacteria_0_collection/mycobacterium_tuberculosis_h37rv/cds/Mycobacterium_tuberculosis_h37rv.ASM19595v2.cds.all.fa.gz

Write a script called `codon_compute.py`. You can download the data outside of the python script or you can include these steps in your script. I already wrote part of this for you in the template code you can start with that executes a `curl` command from within your script.

The code you write will need to process these files in order to print out the following information:

1. The total number of genes in each species.
2. Total length of these gene sequences for each file
3. The G+C percentage for the whole dataset (eg the frequency of G + the frequency of C)
4. Total number codons in each genome.
5. Print out table with three columns: Codon, Frequency in Sp1, Frequency in Sp2