# I529 HW2-Section1

# Chia-Hsuan Chou

## Feb. 19, 2016

## I. Goal

The goal is to test the prediction accuracy of possible *E. coli* genes based on a first order markov chain which generated from 1000 *E coli* genes.

## II. Procedure

First, I collected the entire *E.coli* genome from NCBI and selected 1000 gene sequences from the genome as my training data (attached in the folder which called "gene_1000.fasta"). Next, I used these gene sequences to build a first order markov chain in order to reference it later. Then, I input a test data (attached in the folder which called "500_gene.fasta") in order to test the accuracy of my first order markov chain. In order to test my prediction, I calculated the log-likelihood ratio of $P_i = p_i / p_0$ for my test data. Here notices that I assume every nucleotide is independent to its adjacent codons so that we can multiply every probability of nucleotide based on the markov chain dictionary I made and the outcome will be $p_i$. Also, for the random model of coding DNA, I assume the probability of every nucleotide is $1/4 = 0.25$, so the multiplication of every nucleotide in predicted sequence will be $p_0$. Moreover, I assume that the initial probability for every nucleotide is 0.25. Finally, the result will be the predicted genes and their likelihood and the program gives us two FASTA files with predicted genes and the translated protein respectively.

## III. Result

Figure 1. The output file ("output_gene.fasta") from ECgnfinder_mc.py



Figure 2. The output file ("output_protein.fasta") from ECgnfinder_mc.py

## IV. Discussion

I also calculated the accuracy of my prediction and the accuracy is about 91 percent which I print it on the screen when running the program. Therefore, the accuracy is quite high that we use markov chain model. However, the accuracy is high might because we use all E. coli genes to test the program which means we do not need to consider that it might have non-coding genes inside the test data.