

Nomenclating Aliens

Team Member 1: Omar Sinan (osinan)

Team Member 2: Swapnendu Sanyal (swapnens)

Project Description:

We were recently invaded by aliens from various galaxies. We have managed to somehow fight back but now we want to know more about them. However, we do not have their morphology but only have their DNA samples. We plan on differentiating the species, creating bioweapons tailored to each planet, and retaliate and annihilate them.

The goal of this project is to predict which planet does a particular DNA sequence belong to. We have a multitude of labeled DNA sequences and we want to predict which planet do the invaders belong to.

Project Idea:

The way that the QuAM works is by providing it with a DNA sequence of a particular size as input and using the classifier described below, the QuAM will predict the planet that the DNA sequence belongs to.

We will be building a classifier for this problem using the three methods. Namely, we will be using the normal classifier, k-NN, and k-Means and compare the results to pick the best classifier that solves the problem.

The QuAM will have a simple user interface that allows the user to input the DNA sequence they want to analyze and its size, after running the classifier on the given input, the QuAM will return the prediction of the planet that the DNA sequence belongs to.

Divided Spec

Swapnendu:

- Generate the DNA data by writing a python script.
- Work on one of the classification methods.

Omar:

- Feature extraction from the given data.
- Work on the other two classification methods.

Data Generation

We are going to write a python script that takes the following input:

1. Length of the sequence
2. Number of planets
3. Number of data points

We are going to randomly generate a string of DNA sequence and make it our first cluster. Now, we keep generating clusters that are significantly far from each other for the remaining clusters randomly. Now, we have the clusters, we randomly generate sequences that belong to those clusters. We will output all the DNA sequences with their labels.