

# Interpretable Clustering and Classification of an Imbalanced Dataset of DNA Sequences

Shekh Ahammed Adnan Bashir

Department of Computer Science and Engineering  
Bangladesh University of Engineering and Technology  
1018052026@grad.cse.buet.ac.bd

**Abstract**—Clustering or classification of DNA sequences is more difficult than doing those with any other type of sequences due to various reasons. Such tasks become more difficult when they have to be done on genus or family. Further difficulty is added when the class samples vary in number- that is, the dataset of DNA sequences is imbalanced. However, through using proper feature extraction techniques and dataset sampling technique such tasks are becomes feasible. The dataset we work with in this paper has 9 classes and the number of samples vary from 4 to 1269. In this work, we present a feature extraction technique inspired by popular Natural Language Preprocessing algorithm GloVe [1] to make the classification and clustering of such huge and imbalanced dataset possible. The feature extraction routine is static rather than being a learning one. This eased the interpretation of the machine learning tasks easier. Using interpretable shallow learning techniques, we achieved an accuracy score around 99.8% and a V-measure of 0.5363.

## I. INTRODUCTION

As the rapid development of next generation genome sequencing techniques, newer species are being classified quickly. It is important to know the class (family or genus) that a newly sequenced DNA belong to besides learning how the new species are related to the other ones. Here comes the necessity of DNA sequence classification and clustering. In both classification and clustering, a given collection of items is separated into a few to several subcollections so that the items in one subcollection are as similar as possible with items in two different subcollections are as different as possible. However, they differ in the way they do it. In classification, labeled data is provided to the classifier and the classifier learns an implicit measure of similarity upon seeing the labels of the provided data. This implicitly learned similarity measure can take a very complicated mathematical form. On the other hand, in clustering a measure of similarity is provided explicitly and the clusterer separates the items into subcollections according to that similarity score. The clusterer does not need data labels. The better the similarity measure provided to the clusterer, the better the separation is and hence better the V-measure score is. As the similarity measure used in clustering is already understood, we can gain an insight into the data in the light of that. Broadly speaking, classification discovers the separating boundary in a given collection to divide that into subcollections and clustering is discovering the groups based on a provided similarity measure. Clustering provides us a view of how the items in a collection form homogeneous groups of items and

classification draws decision boundaries among the groups.

A sequence is a list of elements. In a sequence, items of the provided collection are arranged one after another. There might or might not be sequential dependencies among the items. That is, value of an item appearing in the latter positions might or might not depend the value of item or items appearing in the previous positions. This dependency relation varies from problem to problem and this can only either be inferred from the data or provided to the algorithm as parameters. Classification and clustering tasks on the sequences hence are very difficult with data analysis and static algorithms only.

DNA sequences are composed of 4 different nucleotides- Adenine, Cytosine, Guanine, and Thiamine. Therefore, it is very likely that DNA sequence of two species share some common region because of the small state space. There are noises in the DNA sequences. There are intra-class difference of the DNA sequences when it comes classifying genus or families of DNA sequences. Through using preprocessing techniques that maintains the class invariant properties yet transforming all the sequences into sequences of same length or extracting features we can successfully classify or cluster the sequences in such cases.

Interpretability is important to better understand the data so that some static algorithm can be employed to do something of importance with a certain guarantee. When a classifier or a clusterer does their respective tasks, it can do so interpretably or not. Using static feature preprocessing and understanding the parameters and attributes of the machine learning models helps result interpretation. Decision tree based classifiers, probabilistic classifiers, manifold learning algorithms provide interpretable results because they do their in a predefined step by step manner rather optimizing the decision boundary based on different hyperparameters only like black box machine learning models.

In this work, our main focus has been on interpreting the results obtained from the clusterer and the classifier. Specifically, our contributions are:

- Designing a feature extraction procedure that maintains the class invariant property of the DNA sequences.

- Using interpretable machine learning techniques to cluster and classify the DNA sequences.
- Interpretation of the results obtained from the clusterer and the classifier.

## II. PREVIOUS WORKS

Wang et al. proposed two new methods to classify DNA sequences [2]. The first technique relies on the comparison of a given sequence to a group of active motifs discovered from the classes it knows. The second technique generates and matches gapped fingerprints of a given unlabeled sequence with the elements of classes it knows. Stranneheim et al. classified DNA sequences through using Bloom filters to keep track of novelties of the sequence reads[3].

## REFERENCES

- [1] J. Pennington, R. Socher, and C. Manning *Glove: Global Vectors for Word Representation*. EMNLP, 14, page 1532–1543. (2014)
- [2] Wang JT, Rozen S, Shapiro BA, Shasha D, Wang Z, Yin M. *New Techniques for DNA sequence classification*. In J Comput Biol. 1999 Summer;6(2):209-18.
- [3] Henrik Stranneheim<sup>1</sup>, Max Käller, Tobias Allander, Björn Andersson, Lars Arvestad and Joakim Lundeberg. *Classification of DNA sequences using Bloom filters*. Vol. 26 no. 13 2010, pages 15951600, doi:10.1093/bioinformatics/btq230.