# Summary Report: DNA Sequencing Classification Using Machine Learning

Monday, September 30, 2024    7:57 PM

## DNA Sequencing Classification Using Machine Learning

### *Overview:*

This project applied machine learning to classify DNA sequences and predict gene functions. By using k-mer counting to transform DNA sequences into numerical vectors, we trained a **Multinomial Naive Bayes classifier** for accurate gene classification.

### 1. Data Upload and Preparation

The dataset contains labeled DNA sequences, each representing a specific gene function. After uploading the data, we imported necessary libraries, including:

- **pandas** for data manipulation,
- **scikit-learn** for machine learning operations,
- **NumPy** for numerical computations,
- Other libraries like **matplotlib** and **seaborn** for data visualization.

*This project applied machine learning to classify DNA sequences and predict gene functions. By using k-mer counting to transform DNA sequences into numerical vectors, we trained a **Multinomial Naive Bayes classifier** for accurate gene classification.*

### 2. K-mer Counting for Feature Generation

A fundamental concept in this project is converting DNA sequences into a language-like format using **k-mer counting**. K-mers are small, overlapping subsequences of a fixed length, in this case, 6 nucleotides (hexamers).

The function generate_kmers(sequence, k=6) was used to break each DNA sequence into overlapping 6-mers. This process converts varying lengths of DNA into smaller subsequences of uniform length, which are more manageable for machine learning algorithms. These k-mers capture local patterns within the sequences and provide a fixed representation of the data.

#### Why K-mers?

By applying k-mer counting, we ensure that the DNA sequences are represented as numerical vectors. This allows us to apply natural language processing (NLP) tools to the DNA sequences, treating them as if they were sentences composed of "words" (the k-mers). We used **Count Vectorizer** to transform these k-mer sequences into a numerical matrix that the classifier can understand.

### 3. Converting K-mers into Numerical Representation

Once the k-mers were generated, they were converted into a numerical format using **Count Vectorizer** from scikit-learn. Count Vectorizer works by treating each k-mer as a "word" in a document and converting the sequences into a matrix of word counts.

Through parameter tuning, we found that using **n-grams of size 4** yielded the best results, as it allowed the model to capture slightly larger patterns in the DNA sequence data.

### 4. Model Training with Naive Bayes

The processed data was then fed into a **Multinomial Naive Bayes** classifier. This model is particularly well-suited for text-based classification tasks, such as this one where DNA sequences are treated like strings of words.

#### Model Hyperparameters:

- **Alpha ($\alpha$) = 0.1**: This smoothing parameter was optimized for the best performance during training.

#### Why Naive Bayes?

Naive Bayes classifiers are effective in text classification problems, where the assumption of feature independence can lead to efficient and accurate models. In this case, the model treats the presence of each k-mer in a sequence independently, making it a perfect fit for our DNA classification task.

### 5. Model Evaluation

The trained model was evaluated using several key performance metrics:

- **Accuracy**: Measured the overall correctness of the model.
- **Precision**: Indicated the model's ability to correctly classify sequences into the right gene categories without false positives.
- **Recall**: Measured the model's effectiveness at identifying all true positives.
- **F1-score**: Provided a balanced metric between precision and recall.

#### Results:

- The Naive Bayes classifier achieved **98% accuracy** on unseen test data, demonstrating strong generalization and performance.
- The confusion matrix showed very few misclassifications, confirming the model's reliability.

This high performance suggests that the model learned to differentiate gene functions effectively based on the DNA sequences, without overfitting to the training data.

### Conclusion

In this project, I successfully implemented a classification model that predicts gene functions from DNA sequences using machine learning techniques. The use of k-mers and natural language processing allowed us to convert DNA sequences into a format that the Naive Bayes classifier could understand and use for accurate predictions.