

# Protein Stability Prediction (MLB S24) Kaggle Challenge Report

Macrina Lobo

## 1 Data Processing

Exploratory data analysis was performed. This included plotting the distribution of sequence length as well as the histogram of average occurrence of each amino acid in training and test sets. I observed similar distribution on train and test set with each amino acid having good coverage (except U).

I decided to use protein language models (PLMs) for this task. My workflow consisted of tokenization and embedding generation utilizing the last hidden layer of the PLM encoder. Mean pooling was applied at the final layer resulting in a 320 dimension embedding for each input sequence. I tried several *esm2* models[2] and decided to go with *esm2\_t6\_8M\_UR50D* due to having lower number of layers, parameters and embedding dimension thereby leading to faster runtime and lower computational resource requirements. PLM models like ESM2 typically apply padding and truncate the length of the input protein sequences. I truncated my sequence length to a sequence of 1024 amino acids. Alternatively, mean aggregation over every set of 1024 consecutive amino acid sequences in the input instead of truncation was another option. However, the frequency distribution of the first 1024 bases is representative of the entire dataset and consistent on training and test data so truncation might not drastically reduce performance (Figure 1).

Finally, I visualized the embeddings using the first two components of PCA, UMAP computed with 50 principal components, 15 neighbors and a min distance of 0.1. Though there is some separation between high and low stability values on the embeddings, it warrants further downstream transfer learning models.

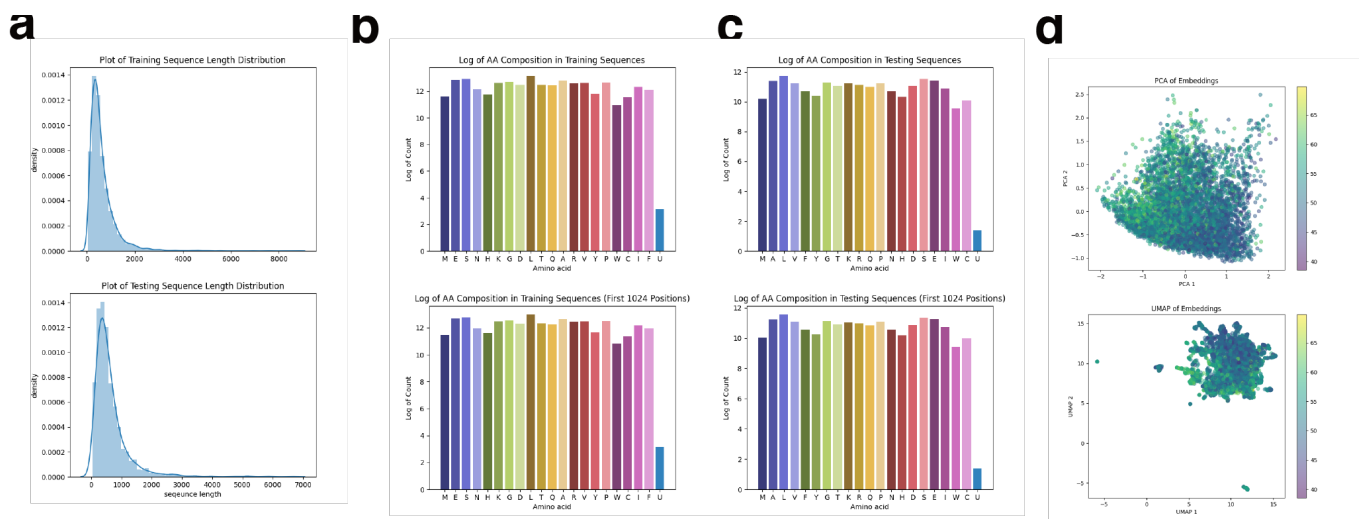


Figure 1: Approaches for Data PreProcessing and Visualization.(a) Histogram showing the sequence length distribution on the training (top) and test (bottom) sets (b-c) Barplot showing the log number of times the amino acid occurs in the whole set (top) and the first 1024 bases which were used here computed on the training set (b) and test set (c)

## 2 Machine Learning Model

From the previous section, there is a need for machine learning models to be trained on the embeddings generated from ESM2. Several downstream transfer learning models were run on the PCA embeddings by tuning the number of principal components. This includes the following:

- KNeighborsRegressor with varying neighbors, weights, algorithm (ball tree, kd tree and brute), and leaf size
- Support Vector Regression with varying C, kernel, degree and gamma.

- Models were run with scikit-learn.

I initially split my training data into 80% training, 20% validation. However, I later decided to use scikit-learn's GridSearch Cross Validation and ran it on the entire training set with the models described in the previous section. I created a pipeline for PCA and the model, defined my grid comprising machine learning models described in the previous section and their tunable parameters and fit the training data and fitness labels with cross validation. I subsequently used the best estimator. All steps were performed locally without gpu on an Apple M2 Max MacBook Pro with memory of 32GB. Use of gpu would lower the runtime.

I tried 60 and 150 principal components. For the other hyperparameters, I tuned them with GridSearch Cross Validation from scikitlearn.model.selection and using a pipeline. I set the evaluation to be Spearmann’s correlation since that was used in the Kaggle evaluation. A total of 108 candidate models alongwith their parameters were tested with 5 fold cross validation totalling 540 fits. A summary of all the tests is available in the all\_results.csv file submitted along with the code. Table 1 contains an overview of the top performing methods for each algorithm tried and their mean performance on the test set during cross validation.

Figure 2: Cross validation results of the top 3 ranked models for each machine learning method applied, namely, KNN-Regression, Support Vector Regression, Random Forest Regression and Multi Layer Perceptron Regression

I achieve high cross-validation score (0.76 spearman correlation) with 60 component PCA and a KNeighborsRegressor using ball tree algorithm, leaf size = 15, 10 neighbors, and weights = distance. I initially got but very poor test correlation (0.02 with the best fit model) but realized this was because my test results were permuted. After fixing this, I got 0.64 test set spearman correlation.

It seems that transfer learning on the ESM2 PLM generated embeddings have generated reasonable test set Spearman Correlation on a simple k-nearest neighbor regression model which even outperformed multi-layer perceptron regression during cross validation. It would be interesting to investigate more advanced models such as VAEs, CNNs, or LSTMs trained on the ESM2 embeddings.

I briefly explored the possibility of using a pre-trained thermostability model EsmTherm[1] on this data. EsmTherm was trained by fine-tuning ESM2 using data from [3]. Applying EsmTherm directly on our Kaggle challenge data yielded poor performance. This is because EsmTherm is trained to predict `delta_g` which is different from the stability measurement in our data. Further, the data in [3] is generated differently from our data. I explored the possibility of fine-tuning *esm2-t6-8M-UR50D* with data from [3] using a LoRA based approach and then using transfer learning on our dataset with these fine-tuned embeddings. However, due to computational and runtime constraints, I could not complete this step though I think it would considerably improve performance.

Code and embeddings have also been uploaded to github at [https://github.com/commacrinalobo2024MLCB\\_thermostability](https://github.com/commacrinalobo2024MLCB_thermostability).

## References

- [1] CHU, K. S., AND SIEGEL, J. B. Protein stability prediction by fine-tuning a protein language model on a mega-scale dataset. *bioRxiv* (2023), 2023–11.
- [2] LIN, Z., AKIN, H., RAO, R., HIE, B., ZHU, Z., LU, W., SMETANIN, N., DOS SANTOS COSTA, A., FAZEL-ZARANDI, M., SERCU, T., CANDIDO, S., ET AL. Language models of protein sequences at the scale of evolution enable accurate structure prediction. *bioRxiv* (2022).
- [3] TSUBOYAMA, K., DAUPARAS, J., CHEN, J., LAINE, E., MOHSENI BEHBAHANI, Y., WEINSTEIN, J. J., MANGAN, N. M., OVCHINNIKOV, S., AND ROCKLIN, G. J. Mega-scale experimental analysis of protein folding stability in biology and design. *Nature* 620, 7973 (2023), 434–444.