

Subject Section

This is a title

Corresponding Author^{1,*}, Co-Author² and Co-Author^{2,*}

¹Department, Institution, City, Post Code, Country and
²Department, Institution, City, Post Code, Country.

*To whom correspondence should be addressed.
Associate Editor: XXXXXXXX
Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation:
Results:
Availability:
Contact: name@bio.com
Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Abstract

The accurate prediction of enzyme commission numbers (EC numbers) is not only crucial for the classification and understanding of newly discovered enzymes but also for completing the annotation of already known enzymes. Therefore, developing a reliable method for predicting EC numbers is of great importance.
However, due to insufficient data, enzyme function prediction using machine learning is an ongoing challenge. In this paper, we propose several methods for predicting enzymes in three different problem categories (Table 1). Throughout the developing of our models, we used a variety of different input features and machine learning algorithms, of which the best will be thoroughly reviewed in this paper.

Table 1. Description of subproblem categories

Level	Description	Best performing method	F1 score
0	Binary classification	Random Forest	score
1	Main class classification	Feedforward neuronal network	score
2	Subclass classification	Feedforward neuronal network	score

catalysis, those reactions would occur too slow or they would even be impossible. Most enzymes are proteins, and each enzyme is assigned a specific function. For example, Oxidoreductases catalyze redox reactions, while Isomerases convert a molecule into its isomer. Consequently, there is a fundamental need for a comprehensive understanding and precise classification of enzymes.
In the past, scientists attempted to categorize enzymes into groups and develop a logical rule set for naming. However, the efforts were hindered by ambiguity. A significant milestone occurred in 1956 with the establishment of an official international commission on enzyme classification (IUBMB). This marked the initiation of the contemporary enzyme classification system that forms the basis for our understanding of enzymes today.
In recent years, computational approaches have emerged as powerful tools for enzyme classification, offering efficient means to navigate through the vast landscape of protein sequences and structures. Enzyme classification, traditionally a labor-intensive task, has witnessed a transformative shift with the advent of computational techniques. This shift is driven by the exponential growth of available biological data, propelled by advancements in high-throughput technologies. Computational methods, particularly machine learning and data-driven models, now stand at the forefront of enzyme classification endeavors, promising to enhance accuracy and efficiency in the annotation of enzymes within large-scale genomic datasets.

2 Introduction

Enzymes serve as crucial biological catalysts, playing an essential role in numerous biochemical reactions within organisms. Without enzyme

3 Methods

The protein data was taken from the UniProt database and used in the CLEAN publication on Enzyme function prediction (see Yu *et al.* (2023)). The mass table of common amino acids was taken from Bio (2021).

3.1 Machine learning algorithms used per level

Our initial goal is to categorize proteins as either enzymes or non-enzymes (level 0). To accomplish this, we’ve opted for the Random Forest, a straightforward yet powerful machine learning method. The choice of Random Forest is driven by its effectiveness in handling classification tasks, making it a well-suited option for our specific protein classification objective. We built a Random Forest Classifier using the scikit-learn library with specific parameters, which will be explained in the training procedure. For ...
For ...

3.2 Data preprocessing in general

Unwanted Sequences Removal We excluded sequences that contain the amino acids "O" and "U" from our dataset.

3.3 Data preprocessing for level 0

As a random forest model relies on multiple decision trees, and each decision tree requires various features, we opted to extract additional information from both the protein sequence and the esm2 embeddings. From the amino acid mass table we computed the mass of protein sequences by adding up the individual masses of their amino acid components. The esm2 embeddings, each represented by a 2560-dimensional vector, underwent statistical analysis. We computed the median, standard deviation, and vector magnitude by aggregating values across all 2560 dimensions for each protein. To simplify the embeddings, we applied Principal Component Analysis (Pedregosa *et al.* (2011)) separately to enzyme and non-enzyme datasets, retaining 90% of the variance. This process yielded reduced dimensions of 397 for enzymes and 369 for non-enzymes. Therefore, we reduced the dimensions to 397, providing a streamlined representation of the protein embeddings while retaining crucial information for both enzyme (see figure 1) and non-enzyme (see figure 2) datasets.

By combining the information from the proteins’ sequences, masses, and embeddings, we created a Pandas DataFrame that concatenates enzymes and non-enzymes, including 401 features: mass, embeddings median, embeddings standard deviation (std), embeddings magnitude, and dimension 1 to 397 of the reduced embeddings.

We also used the method SelectFromModel from the scikit-learn library (see Pedregosa *et al.* (2011)) to select features based on importance weights. The refined set of input features contains 'Mass,' 'Emb median,' 'Emb std,' 'Emb magnitude,' and 'PCA 1' through 'PCA 47' (see figure 3) excluding 'PCA 29', 'PCA 31', 'PCA 35', 'PCA 36', 'PCA 39' to 'PCA 42', 'PCA 44' and 'PCA 46' (where 'Emb' represents embeddings and 'PCA X' signifies the X-th dimension of the reduced embeddings).

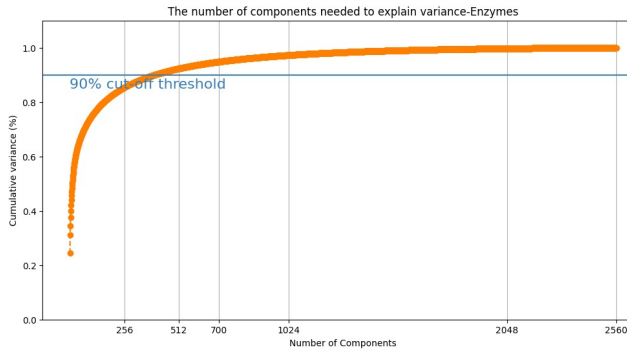


Fig. 1. The number of components needed to explain the variance in the enzyme dataset

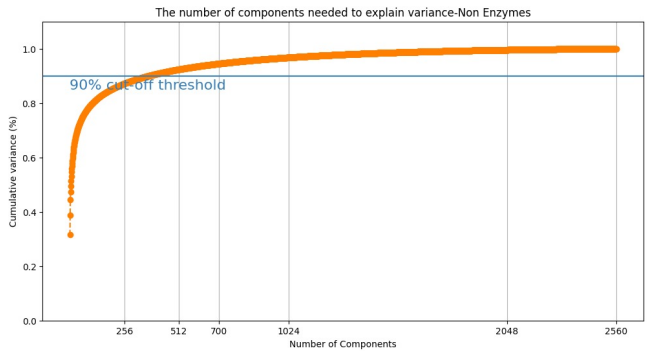


Fig. 2. The number of components needed to explain the variance in the non enzyme dataset

Mass	Emb median	Emb std	Emb magnitude	PCA 1	PCA 2	PCA 3	PCA 4	PCA 5	PCA 6	...
60153.711	-0.002189	0.227337	11.502516	0.450330	0.056761	-0.152714	-0.562537	-0.342939	-0.214768	...
81547.542	-0.002620	0.240143	12.150464	-0.087801	-0.501041	0.027421	-0.060358	-0.472940	-0.169462	...
73156.055	-0.001840	0.243909	12.341031	-0.490790	0.191398	-0.227236	-0.447605	0.192044	-0.332910	...
16788.450	-0.003522	0.263979	13.356579	0.409453	-0.436421	0.680466	0.608248	-0.913846	-0.340269	...
17145.211	-0.004006	0.273501	13.838382	-1.157048	-0.699554	-1.003535	0.209268	-0.485984	0.239973	...
...
38569.516	-0.002372	0.237258	12.004517	1.250809	0.068282	-0.679364	-0.326445	-0.044144	0.019516	...
62801.437	-0.003409	0.240624	12.174803	1.049215	-0.032362	-0.520308	0.227688	-0.107257	0.025082	...
42719.902	-0.003696	0.260407	13.175781	-0.141801	-0.123650	-0.129375	0.378490	-0.209925	-0.149825	...
39507.175	-0.004268	0.227770	11.524400	1.324474	-0.121695	-0.079250	0.418808	0.021615	-0.275152	...
90720.702	-0.004878	0.280515	14.193187	-1.356220	-0.648414	-0.067346	0.474130	0.131560	0.340424	...

Fig. 3. The number of components needed to explain the variance in the non enzyme dataset

3.4 Training procedure

For all our models we divided the data table into two parts: a training set and a validation set, using a random state of 42 for consistency. The training set contains 70% of the original data, while the validation set holds the remaining 30%. This separation allows us to train our models on a subset of the data and then assess its performance on a different subset to ensure its generalization to new, unseen data.

3.4.1 Random Forest

For the Random Forest classifier we addressed the imbalance between the number of non-enzymes and enzymes, where non-enzymes outnumber enzymes approximately fourfold, by duplicating the enzyme data in the training set four times. This duplication ensures a more balanced dataset, allowing the model to be trained on an equal representation of both classes.

The classifier consists of a total of 200 decision trees. Each tree has a maximum depth of 16, and a node is designated as a leaf only if it has a minimum of 8 samples. The random state parameter is set to 42, ensuring reproducibility in the model’s construction.

3.4.2 Level 1 FNN
3.4.3 Level 2 FNN
3.5 Validation on test dataset
3.5.1 Scoring metrics

precision = TP / (TP + FP) (1)

recall = TP / (TP + FN) (2)

f1-score = 2 * (precision * recall) / (precision + recall) (3)

weighted f1-score = sum_{i=1}^N (Number of samples in class i / Total number of samples) * f1-score_i (4)

MCC score = (TP * TN - FP * FN) / sqrt((TP + FP) * (TP + FN) * (TN + FP) * (TN + FN)) (5)

The MCC score serves as an indicator for the accuracy of binary classifications, with values ranging from -1 to 1. A score of 1 signifies a flawless prediction, -1 indicates a completely inaccurate prediction, and 0 suggests predictions at random.

To assess the performance of our models, we emphasize the F1 score and MCC score. Given the imbalance in the test set, we also considered the weighted F1 score to ensure a more equitable evaluation of the model's effectiveness.

4 Results

4.1 Random Forest Level 0

Using the Random Forest model on the "new" test set, we achieved accurate predictions for about 90% of positive cases (enzymes) and 98% of negative cases (non-enzymes), even with an imbalanced test set comprising 392 enzymes and 9876 non-enzymes.

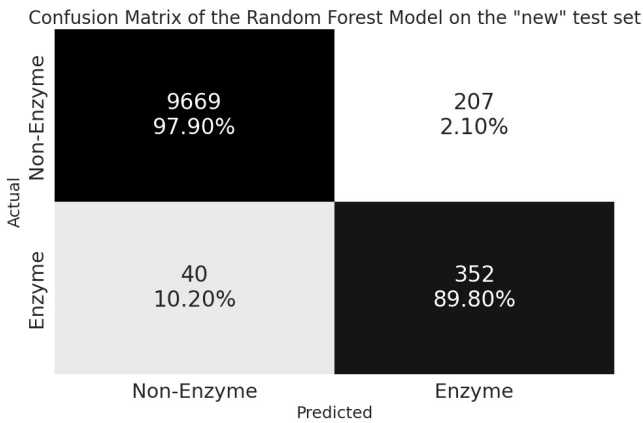


Fig. 4. Confusion Matrix of the Random Forest Model on the "new" dataset

Additionally, we attained an Accuracy of 97.6%, a weighted f1-score of 97.8%, and an MCC-score of 74.1%. Figure 5 illustrates that our Random Forest model significantly outperformed the random baseline.

Furthermore, when we tested our model on a distinct test set, the "price" dataset, it demonstrated a strong performance: only one enzyme was misclassified (see figure 6). This reinforces the robustness and generalizability of our Random Forest model across diverse datasets.

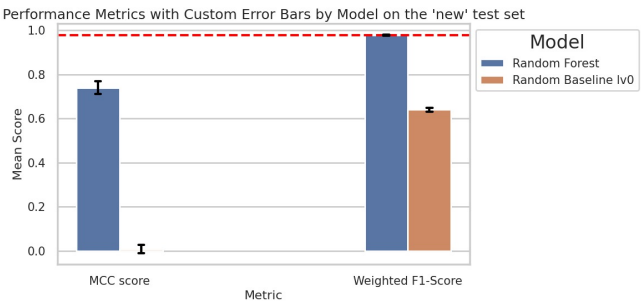


Fig. 5. Level 0 model comparison of Random Forest and baseline on "new" dataset

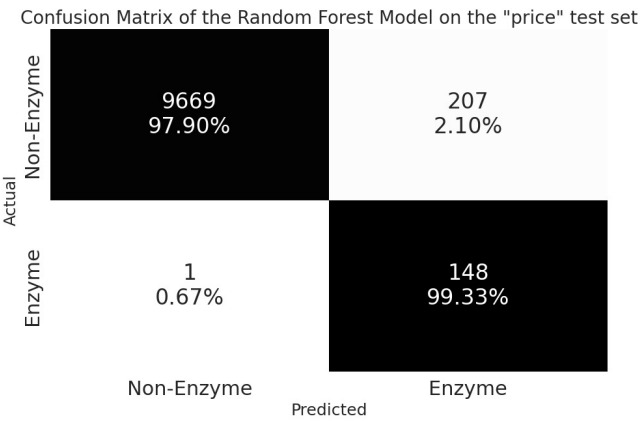


Fig. 6. Confusion Matrix of the Random Forest Model on the "price" dataset

4.1.1 Support Vector Machine
4.2 Level 1 performance
4.3 Level 2 performance
5 Discussion
6 Conclusion
7 Supplementary Information

7.1 K-nearest neighbors algorithm using ncd vectors

A less popular approach of transforming string like input features into numerical values is the normalized compression distance (ncd) algorithm. The ncd algorithm is based on the idea that the similarity of two strings can be measured by the amount of information needed to describe one string given the other string. The ncd of two strings x and y is defined as follows:

ncd(x, y) = (C(xy) - min(C(x), C(y))) / max(C(x), C(y)) (6)

where $C(x)$ is the length of the compressed string x and $C(xy)$ is the length of the concatenated string xy .

We implemented this algorithm in python using `gzip`, which is a loss less compression algorithm based on a combination of LZ77 and Huffman encoding. Rigler *et al.* (2007)

The ncd algorithm was used to transform the amino acid sequences into numerical vectors by comparing each sequence to all other sequences in the training dataset. This resulted in a n -dimensional numerical vector for each sequence, where n is the amount of sequences in the training dataset where each position in the vector represents the ncd of the sequence to the corresponding sequence in the training dataset. These vectors were then used as input for the k-nearest neighbors algorithm. Due to the exponential

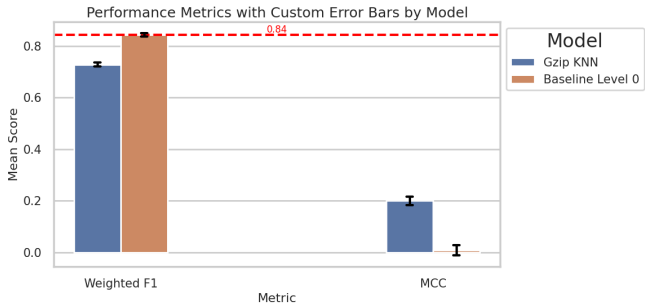


Fig. 7. Performance on test dataset compared to random baseline

computational complexity of the ncd algorithm, we had to under sample the non enzyme dataset to match the amount of samples in our enzyme dataset, meaning the positiv in the train dataset were balanced.

When inferring unseen data, the ncd input vector was calculated by comparing it to all sequences in the training data set, thus also resulting in a n -dimensional numerical vector. This means that the performance on new data is largely dependent on the training data set.

The performance of the k-nearest neighbors algorithm using ncd vectors compared to a random baseline is shown in figure 7. Although the mean F1 score of the k-nearest neighbors algorithm using ncd vectors lies

at 0.728, it did not perform better than the random baseline, which had a F1 score of 0.843. This indicates that the ncd approach has a worse precision and recall than the random baseline. At the same time both classifiers have a low MCC score of 0.2 and 0.01 respectively, which indicates that both classifiers are not better than random guessing.

The reason for the poor performance of the k-nearest neighbors algorithm using ncd vectors is most likely due to the ncd algorithm not being suited for protein sequences as shown in Matsumoto *et al.* (2000) as well as the test dataset not being balanced, while the training dataset was.

References

(2021).

Matsumoto, T., Sadakane, K., and Imai, H. (2000). Biological sequence compression algorithms. *Genome Informatics*, **11**, 50–51.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, **12**, 2825–2830.

Rigler, S., Bishop, W., and Kennings, A. (2007). Fpga-based lossless data compression using huffman and lz77 algorithms. In *2007 Canadian Conference on Electrical and Computer Engineering*, pages 1235–1238.

Yu, T., Cui, H., Li, J. C., Luo, Y., Jiang, G., and Zhao, H. (2023). Enzyme function prediction using contrastive learning. *Science*, **379**(6639), 1358–1363.