

Applications of Transformer Networks in Bio- and Cheminformatics

Worksheet 1

Submission Deadline:	22. April 2025, 23:59
Discussion of solutions:	29. April 2024, 12:30 - 14:30

Submission Instructions:

- Create a single Jupyter notebook for Exercises 1, 2, and 3. This Jupyter notebook should contain your code and the results for these exercises.
- If you want to submit multiple files, zip them together. The uploaded file should have the following filename “lastname_studentID_worksheet1”.
- Each submission must include an answer to Exercise 1.4. You can write the answer in your Jupyter notebook or in a text file.
- Submit your solutions by uploading the Jupyter notebook to <https://uni-duesseldorf.sciebo.de/s/hCt1rTP23EeWmUC>.

Exercise 1.1 *Data Preprocessing* (20 Points)

Each enzyme has a preferred pH at which it does its job most efficiently. We refer to the ideal pH as the enzyme’s optimal pH. Above or below this pH, the activity of the enzyme may decrease, making it less effective. In this worksheet, we want to predict experimentally obtained optimal pH values.

- (a) Download the datasets *phopt_training.fasta*, *phopt_validation.fasta*, and *phopt_testing.fasta* from <https://zenodo.org/records/8011249>. Load the datasets into Python: Create three Pandas DataFrames *train_df*, *val_df* and *test_df*, each with two columns “Sequence” and “optimal pH”. The “Sequence” column should contain the protein amino acid sequences as strings; the “optimal pH” column should contain the optimal pH as a float value.
- (b) Plot histograms showing the distributions of sequence length and optimal pH in the training set.

Exercise 1.2 *Fitting and validating prediction models*

(40 Points)

- (a) Define a function to list all possible k-mers for a given k for the amino acids 'A', 'C', 'D', 'E', 'F', 'G', 'H', 'I', 'K', 'L', 'M', 'N', 'P', 'Q', 'R', 'S', 'T', 'V', 'W', 'Y'. The function should return all k-mer strings as a list.
- (b) Given a list of k-mers, define a function that, for a protein sequence, returns an array containing the number of occurrences of each k-mer. The frequencies should be listed in the same order as the k-mers are stored in the list given as input.
- (c) Calculate all 1-mers for the proteins in *train_df*, *val_df* and *test_df* from Exercise 1.1. Convert the absolute occurrences of these 1-mers to frequencies, i.e., divide by the total number of 1-mers in a sequence. Store the resulting vectors in a new column column "1-mer" for all dataframes.
- (d) Repeat (c) for 2-mers and 3-mers.
- (e) Fit linear regression models on the training set using k-mer frequencies as model inputs (k=1,2,3) to predict the optimal pH of a protein. Validate the models on the validation set using the mean squared error and the coefficient of determination R^2 . Report the best model based on the performance on the validation set.

Exercise 1.3 *Using Pfam domains*

(40 Points)

I generated a Python dictionary that maps protein sequences from the pH dataset to their Pfam domains. You can download the dictionary from https://github.com/AlexanderKroll/DL4Molecules_Course_2025/data/Worksheet1, and you can open it with:

```
seq_to_pfam = np.load('../data/pH/seq_to_pfam.npy',  
    allow_pickle=True).item()
```

- (a) Map each sequence to a binary vector that encodes which domains are present in the sequence. Only encode domains that occur at least twice in the entire data set.
- (b) Train a linear regression model with L2 regularization coefficient $\alpha = 1$:

```
from sklearn.linear_model import Ridge  
model = Ridge(alpha=1.0)
```

- (c) Validate the model on the validation and the test set.

Exercise 1.4 *Use of LLMs*

(0 Points)

State for which exercise you have used LLMs (large language models) such as ChatGPT or GitHub Copilot. State which tools you have used and for which steps. This answer does not influence how many points you receive for your submission.