

Machine Learning 2024 Course - Homework

Francesco Trovò, Vanja Miskovic, Stefano Samele
{francesco1.trovo, vanja.miskovic, stefano.samele}@polimi.it

April 29, 2024

The current homework is focused on the prediction of *taxonomic identity and genetic composition based on codon usage bias levels*. In what follows we detail the context and the task required to earn grades from the homework. The data have been taken from *Khomtchouk, Bohdan B. "Codon usage bias levels predict taxonomic identity and genetic composition." bioRxiv (2020)*.

1 Background

The coding DNA of a genome describes the proteins of the organism in terms of 64 different codons that map to 21 different amino acids and a stop signal. Different organisms differ not only in the aminoacid sequences of their proteins, but also in the extents in which they use the synonymous codons for different amino acids. The inherent redundancy of the genetic code allows the same amino acid to be specified by one to five different codons so that there are, in principle, many different nucleic acids to describe the primary structure of a given protein. Coding DNA sequences can carry information beyond that needed for encoding amino acid sequence. Thus, one may ask: is it possible to classify some properties of nucleic acids from the usages of different synonymous codons rather than, with much greater computational effort, from individual nucleotide sequences themselves?

This data set enables a preliminary analysis on this topic. In particular, codon usage frequencies from several organisms are studied to identify if they can be used to classify codon usage:

1. in terms of viral, phageal, bacterial, archaeal, and eukaryotic lineage;
2. classifying codon usage by cellular compartment (nuclear, mitochondrial, and chloroplast DNA in their respective organelles).

2 Dataset Composition

The data set consists of a total of 13028 organisms (samples) classified both based on the organism species (11 classes) and on the DNA type (11 classes).

Features are all real valued and no data is missing, except from the AGA codon frequency on test samples.

Features Every sample is described by 67 attributes:

- SpeciesName : the species of the organism;
- SpeciesID : an identifier of the species;
- Ncodons : the total number of codons in the entry;
- UUU - AUG : the codon frequencies (the codon occurrence count divided by the total number of codons Ncodons in the entry).

Classes The dataset features two classification problems. In the first, the goal is to classify samples in 11 different classes based on the value of the attribute Kingdom (11 classes): a three letter code referring to the organism species. In particular, identifiers are: 'arc' (archaea), 'bct' (bacteria), 'phg' (bacteriophage), 'plm' (plasmid), 'pln' (plant), 'inv' (invertebrate), 'vrt' (vertebrate), 'mam' (mammal), 'rod' (rodent), 'pri' (primate), and 'vrl' (virus).

The second task is similar, but samples are classified based on the DNAtype (11 classes): an identifier referring to the type of DNA sample: 0 (genomic), 1 (mitochondrial), 2 (chloroplast), 3 (cyanobacteria), 4 (plastid), 5 (nucleomorph), 6 (secondary endosymbiont), 7 (chromoplast), 8 (leucoplast), 9 (NA), 10 (proplastid), 11 (apicoplast), and 12 (kinetoplast).

Train/Test splits The dataset is provided already split into train (10422 samples) and test (2606 samples) sets:

- train.csv : 10422 \times 67 dataset;
- test.csv : 2606 \times 66 dataset.

While the training set comes with all the 67 codon frequencies for each sample, a malfunctioning in the analysis procedure caused the measurements of the AGA codon frequency to be lost in all test samples. Features are all real valued and no data is missing, except from the AGA frequency on test samples.

Requests

1. Perform a preliminary analysis on the data. For instance, but not limited to, visualize samples, identify if features (i.e., codon frequencies) are correlated, determine which are most correlated with the target classes, and inspect the distribution of samples among classes. Using clustering, study if there are structures in the data that allow samples from different classes (both DNA type and Kingdom) to be easily identified. Compare the performance of different clustering algorithms and distance measures using the metrics presented during the course.

2. Perform classification to classify organisms into the 11 Kingdom classes. Perform features selection, compare different algorithms and identify the one that works the best on this dataset. Finally, test the performance of the best algorithm on the provided test set.
3. We want to recover from the data loss of the AGA frequency on test samples. Train a regressor which can predict the value of the AGA feature given the remaining ones. Compare different regression algorithms for this task. Since AGA features are missing in test samples, use only the training data for this step and make use of robust evaluation techniques to compare algorithms.
4. Use the regression model trained at the previous step to recover the AGA codon frequency, by predicting the AGA value of each test sample. Determine if the test performance of the best model found at step (2.) improves if the AGA frequency values are also used for prediction.

Note: As the data set contains several samples, some algorithms and visualizations may be particularly slow to run. You can opt for selecting a smaller portion of the data to perform part of the analysis (part of the preliminary analysis may be devoted to select a meaningful portion on the data set), or other techniques that you think could be useful to mitigate the issue.

Submission The final report should be sent by mail (to both the TAs and the lecturer) before the student takes the written exam. You should:

- Create a Jupyter notebook to answer all the requests, using the libraries presented during the laboratory classes;
- Include Name, Surname and Student ID in the notebook;
- You are free to use the structure that you prefer within the notebook. However, please use markdown cells (Cell > Cell Type > Markdown) to insert section titles and clearly identify the different requests. You are free to add subsections to make the notebook more readable;
- Add text cells (markdown) to briefly explain what you did and why, and to help you answer the requests;
- Please, check that the notebook can execute correctly before submitting your work. Press Kernel > Restart & Run All and check that all cells execute correctly without errors;
- The output of notebook cells should be included in the submitted notebook.

Evaluation The project evaluation is not performed based on the (clustering, classification, regression) scores you obtain in each request, but rather on the soundness of your analysis and the choices you made to solve issues (if any) you may encounter. A good project follows the evaluations discussed during the lectures and lab sessions, while also adapting and integrating the analysis based on the project. The following is the evaluation criterion used in the previous year's project. Notice that it may slightly vary in this year evaluation, but you can use it as a rough reference.

Points for each task: Clustering - Classification - Regression

- 0: Completely wrong, or one key concept is not taken into account;
- 1: Ok, but poorly discussed, e.g., just a copy-paste of labs with no comments, or some conceptual errors;
- 2 Well done, good comments and analysis, adapted code/analysis, he/she interpreted the results.

Penalties

- 0: Does compile with no error;
- -0.5: Does not compile, but it is an easy fix to make it work;
- -1 Does not compile, and it would require more than 10 mins of debugging the errors