

Bioinformatics

Protein Classification By feature extraction

Sofiane MAHIOU

Computer Science, UCL, London, WC1E 6BT, UK

Abstract

Aim : The goal of this assignment is to provide an automated system that is able to classify proteins (Amino Acid sequences) into four classes each being a subcellular locations : [**Cytosolic, Secreted, Nuclear, Mitochondrial**]

Results : Using a *Random Forest Classifier* we manage to reach a **67% cross-validation accuracy**.

Improvements: In order to improve the results of the classifier, deepening the feature extraction method seems to be the way to go. Another method would be to use neural network techniques.

Contact: ucabsm1@ucl.ac.uk

1 Introduction

Currently, There is a growing need for fully automated methods to analyse amino acids sequences. One of the process that need to be automated is **the identification of the protein's subcellular location**. This problem can be splitted into two sub problems :

- **feature extraction:** the goal of this task is to choose the features that would allow an efficient classification, to be more precise, the chosen features should allow to easily separate the sequences into classes or groups which will then be matched with the various subcellular locations
- **classification:** once the features obtained, it is then necessary to choose a fitting classification algorithm that will use the various features selected as a *vector representation* of each sequence that will then be fed to the classification algorithm during both training and testing.

It is however, possible to avoid splitting the problem into two sub-problems by using methods that have been designed to classify sequences of variables lengths such as :

- **HMM**
- **Recurrent neural networks & Seq2Seq Models**
- **1D Convolutional Neural Networks**

Although these methods usually yields better results than the methods presented before, the results obtained are far harder to interpret as these systems behave as "**black boxes**" and it's quite difficult to interpret what was *learned*.

Therefore, The first approach was used in order to ease the analysis of the results and the task was splitted into a **feature-extraction** task and **classification** task. Therefore, in order to indentify what features might be useful to this problem, a research phase was realized where several research papers on the same subject have been studied, and used as a reference to select several features.

2 Sequence preprocessing :

It seems important to mention that before proceeding the **feature extraction** phase, it was necessary to preprocess the data due to the presence of unexpected characters : **U, B, X** in the amino acid sequences.

- **X:** Given that "X" refers to "any amino acid" we randomly replace it by a given amino acid among the 20 amino acids
- **B:** "B" refers to either *asparagine* and *aspartic acid*. It is therefore automatically replaced by one or the other
- **U:** the "U" amino acid is simply removed from the sequence for lack of a better solution

This preprocessing step, will only affect 64 sequences out of more than 9000, therefore these changes are unlikely to heavily influence the results but will allow and easier implementation of the various features.

- **Discriminative Descriptors:** this is another method that was presented in Saidi *et al.* (2010), it behaves similarly to tf-idf but uses the additional information that several classes exist. [_____ TO COMPLETE _____]

3.6 Nuclear Export Signals: **Used**

- This feature that is presented in Xua *et al.* (2012), it describes the following pattern as an efficient discriminative pattern : $\phi_1 - X_3 - \phi_2 - X_2 - \phi_3 - X - \phi_4$. Positions ϕ_3 and ϕ_4 of this prevalent pattern are dominated by the five traditional hydrophobic residues **Leu[L]**, **Ile[I]**, **Val[V]**, **Met[M]**, and **Phe[F]**.
- This feature did lead to improvements, however they were as significant as expected.

3.7 Nuclear Localization Signals: Used

This feature has also been presented in Xua *et al.* (2012), it refers to the count of subsequences of at least 5 **positively charged amino acids** .ie meaning one of the following : **lys[K], arg[R], his[H]**. Again this feature lead to slight improvement over all models.

3.8 Protein's properties: **Used**

In addition to what was presented above the following properties were attempted as suggested by Q.-B. Gao *et al.* (2005)

- **Hydrophobicity**
- **Aromaticity**
- **Molecular Weight**

3.9 Begning and End of sequences:

Each feature presented above was computed for the **full sequence** as well as the **first 50 amino acids** and the **last 50 amino acids** of each sequence. This aims to identify trends and patterns not only overall but also specific to the beginning and the end of the sequences. Indeed, the length of sequences would make it difficult to extract information solely related to the beginning of the sequence therefore isolating the most probably relevant subsequences seems to be a proper way of removing the "noise" due to then length of most amino acid sequences.

4 Classification Methods

Several classification methods are available, using previous knowledge about **Machine Learning - Classification** problems as well as classifiers referenced by other research papers of this same field. This lead to try the following classifiers :

- Logistic Regression Classifier
- Random Forest Classifier
- SVM Classifier
- Ridge Regression Classifier

Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text. Figure 2 shows that the above method
Text Text Text Text Text Text Text Text Text Text Text Text Text. ? might want
to know about text text text text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text Text Text Text Text Text Text Text. Figure 2
shows that the above method Text Text Text Text Text Text Text Text Text
Text Text Text. ? might want to know about text text text text Text Text
Text Text Text Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text. Figure 2 shows that the above method Text Text Text

Text Text Text Text Text Text Text Text Text Text Text Text Text
Text Text Text Text Text Text Text. Figure 2 shows that the above method
Text Text Text Text

Text Text Text Text Text Text Text. ? might want to know about text
text text text

This work has been supported by the... Text Text Text Text.

R. Saidi, M. Maddouri and EM. Nguifo (2010)
Protein sequences classification by means of feature extraction with substitution matrices, *BMC Bioinformatics*

Q-B Gao, Z-Z Wang, C Yan, Y-H Du (2005)
Prediction of protein subcellular location using a combined feature of sequence, *FEBS Letters*

D Xua, A Farmer, G Colletta, N V. Grishin, Y M Chooka (2012)
Sequence and structural analyses of nuclear export signals in the NESdb database, *MBoC*