

# Prediction of Signal Peptide Cleavage Site Using Supervised Learning

Contact: Sarah ([sarah.berkemer@polytechnique.edu](mailto:sarah.berkemer@polytechnique.edu))

Conceptual design by Philippe Chassignet

## 1 Introduction

Protein targeting is the mechanism by which some proteins are directed to their appropriate destinations in the cell or outside. The delivery process is carried out, based on information contained in the protein itself. This information is often identified as a short sequence of amino acids located near the N-terminal side of the protein. This sequence, known as the *signal peptide*, extends more or less toward the C-terminal side of the protein, up to a position named the cleavage site. Once a protein had been properly directed, its N-terminal part is generally cleaved out at this position, while the C-terminal part leads to the protein as observed in its mature form. Correct targeting of proteins is crucial for life and the identification of signal peptides and their cleavage site is of interest to the design of new drugs.

## 2 Assignment

We consider the problem of predicting the position of cleavage site in proteins, relying on patterns learned from a training set. Here a protein is known as its primary sequence, i.e. the sequence of its amino acids, which is given starting from the N-terminal side. Each amino acid is denoted by a letter code and, as there are 20 standard amino acids, plus some peculiarities, most of upper case letters are used in this encoding. Let  $\mathcal{A} = \{A, \dots, Z\}$  be this alphabet.

For instance, the following sequence is the beginning of a protein, where the cleavage site is marked as the bond between the two underlined AR amino acids :

MASKATLLLAFTLLFATCIARHQQRQQQNQCQLQNIEALEPIEVIQAEA...

Then, for the purpose of this project, one will simply work on such sequences of letters. One has now to program some learning algorithms, to tune them, and to evaluate their performance. A simple statistical model, will be used first to establish a reference. One will then try to improve accuracy with some specific kernel functions for Support Vector Machines.

## 3 Localization of a Cleavage Site as a Classification Problem

It has been shown in [1] that the cleavage site may be characterized by amino acids in its close neighborhood. Then, this neighborhood is defined as the  $p$  amino acids before and the  $q$  after. Note that the cleavage site is not an amino acid but the bond between two consecutive amino acids. So, values  $p$  and  $q$  define a neighborhood of  $p + q$  amino acids in length. The work in reference suggests up to  $p = 13$  and  $q = 2$ , but different values may also be experimented as meta-parameters of the problem.

We are facing a binary classification problem. Given any whole protein sequence  $(a_i)_{i=0,\dots,\ell-1}$ , and any position  $j$ , where  $p \leq j \leq \ell - q$ , the word  $a_{j-p}a_{j-p+1} \cdots a_{j-1}a_j \cdots a_{j+q-1} \in \mathcal{A}^{p+q}$  should be enough to decide whether the bond at position  $j$ , between  $a_{j-1}$  and  $a_j$ , is a cleavage site or not.

In some cases, one may get several positive answers along a same sequence. This problem can be considered as an extension to the project

### 3.1 Simple Statistical Model Using a Position-Specific Scoring Matrix

Over a set of  $N$  sequences, with a known cleavage site for each, one can first count  $c(a, i)$ , the number of occurrences of each amino acid  $a \in \mathcal{A}$ , at every position  $i \in \{-p, \dots, q-1\}$ , relative to the corresponding cleavage site. Then, for each  $a$  and  $i$ , let define  $f(a, i) = c(a, i)/N$ , the observed frequency of amino acid  $a$  at the relative position  $i$ .

In a same way, by counting over the whole length of given sequences, one can compute the observed general background frequency  $g(a)$  of amino acid  $a$  in the given set, regardless of the position. However, it must be noticed that the very first amino acid at the beginning of a sequence is almost always an M, because it corresponds to the transcription of the start codon. Also, one will not count letters on this first position to avoid a bias.

These frequencies will be used as estimated probabilities to compute the probability of a given word to be located at a cleavage site, under an independent model. We rather use the logarithm of probabilities to go on additive calculations.

Then,  $\forall a \in \mathcal{A}, \forall i \in \{-p, \dots, q-1\}$ , we define  $s(a, i) = \log(f(a, i)) - \log(g(a))$ . Also, as zero counts may occur, pseudo-counts [3] must be used. Finally, for any word  $w = a_0a_1 \cdots a_{p+q-1}$ , the  $q-1$  score defined as  $\sum_{i=-p}^{q-1} s(a_{p+i}, i)$  may tell whether  $w$  is the neighborhood of a cleavage site or not. A simple thresholding (to be tuned) is then enough to define a simple binary classifier.

### 3.2 Some SVM Kernels

One must implement and experiment some classification algorithms with SVM around the following ideas. It should be easy to achieve better results than with the previous statistical model.

Kernels for SVM are basically positive (semi-)definite functions  $\mathbb{R}^d \times \mathbb{R} \rightarrow \mathbb{R}$ , used as a generalization of the inner product for space dimension  $d$ . Instead, we need some kernel function  $K : \mathcal{A}^{p+q} \times \mathcal{A}^{p+q} \rightarrow \mathbb{R}$ , defined on words of length  $p+q$ . First solution is to have an explicit mapping  $\Phi : \mathcal{A}^{p+q} \rightarrow \mathbb{R}^n$ , which allows to define a kernel as  $\forall u, v \in \mathcal{A}^{p+q} : K(u, v) = \Phi(u) \cdot \Phi(v)$ , where  $\cdot$  is the inner product on  $\mathbb{R}^d$ . As  $d$  may be large, an interesting property of some kernels is that one does not need to explicitly calculate  $\Phi$  and, then, the inner product.

A direct approach for using SVM on fixed length words, consists in encoding each word as a point in a Euclidean space. A simple way to encode a single letter is to form a vector of length  $26n$ , containing only zeros, but a one at the  $i^{\text{th}}$  position when encoding the  $i^{\text{th}}$  letter of the alphabet. Then the encoding of every word of  $n$  letters is the vector of length  $26n$ , obtained by concatenating the  $n$  vectors for the letters of this word. Such an encoding then allows a straight linear SVM, or one of the popular polynomials and Gaussian RBF kernels.

One can observe that the dot product of two vectors, as defined by the previous encoding, simply counts the number of common letters between the two corresponding words. This leads to a dedicated implementation with a kernel defined as a similarity function of two words. Then, one should retrieve exactly the same results, as for the extensive encoding.

For a more flexible criterion, one can use the similarity defined by a substitution matrix [4]. Substitution matrices are built on the observed mutations of amino acids within proteins having

a similar function. This yields a matrix of scores of similarity  $M(x, y)$  between any pair  $(x, y)$  of amino acids. Then, to get a similarity score  $S(a, b)$  between two words,  $a = a_0 \cdots a_{n-1}$  and  $b = b_0 \cdots b_{n-1}$ , one may define  $S(a, b) = \sum_{i=0}^{n-1} M(a_i, b_i)$ , the sum of the corresponding letter scores along the two words. It is also a simplification of a more general alignment algorithm, but where no insertion/deletion is considered. Note that such a score is not related to a distance and that it does not satisfy Mercer's conditions. But it is safe to use it as a pseudo dot product to define a RBF kernel, for instance.

As another approach, one can use a probabilistic kernel, based on the corrected frequencies  $s(a, i)$  [2]. In short, for two words,  $a = a_0 \cdots a_{p+q-1}$  and  $b = b_0 \cdots b_{p+q-1}$ , the logarithm of such a kernel can be defined as :

$$\log K(a, b) = \sum_{i=-p}^{q-1} \phi_i(a_{p+i}, b_{p+i})$$

where

$$\phi_i(x, y) = \begin{cases} s(x, i) + s(y, i) & \text{if } x \neq y \\ s(x, i) + \log(1 + e^{s(x, i)}) & \text{if } x = y \end{cases}.$$

At last, since the two kernels defined above are based on very different properties, one can also try to combine them (as a weighted sum or as a product) into a single kernel.

## 4 Expected Work

This project requires learning from a given sample set and then predicting whether any given word is the location of a cleavage site or not. For that, you have to write program (or many), which allows the choice of a kernel and the tuning of meta-parameters, as the size of the neighborhood around the cleavage site, the  $C$  parameter of SVM and any other parameter, for example the parameter of a RBF kernel.

You have to use an existing SVM library, but implementing your own kernels. You will see the usage of C++ `libsvm` library during the exercise sessions, but a large part of the work can be started before. Note that `libsvm` also exists for Java and `scikit-learn` for Python also allows the definition of custom kernels. Then, these three languages are possible (and allowed) for this project. Execution time is important as tuning a model will require many consecutive runs.

One can then slide a window onto a given sequence, examining in turn each position as a potential cleavage site, to locate the real one, or all those classified as positive, as you can possibly locate more than one for some sequence. So, another program must implement your *best* classifier, already tuned, take a series of sequences as an input file, and output the predicted positions of cleavage site.

## 5 Data Sets

The BLOSUM62 matrix can be downloaded from <https://www.ncbi.nlm.nih.gov/Class/FieldGuide/BLOSUM62.txt>

The project folder contains three files containing sequences from various organisms, where the signal peptide is followed by a part of the amino acids of the mature protein:

- `EUKSIG_13.red` contains 1005 sequences from eucaryotes,
- `GRAM+SIG_13.red` contains 140 sequences from Gram-positive prokaryotes.

- GRAM-SIG\_13.red contains 263 sequences from Gram-negative prokaryotes.
- SIG\_13.red is a concatenation of the three files listed above.

Each entry in these files consists in exactly three lines. This is an example of a protein description:

55 2SS8\_HELAN 25 ALBUMIN 8 PRECURSOR (METHIONINE-RICH 2S PROTEIN) (SFA8).  
MARFSIVFAAAGVLLLAMAPVSEASTTTIITTIIEENPYGRGRTESGCYQMEE  
SSSSSSSSSSSSSSSSSSSSSSSCMMMMMMMMMMMMMMMMMMMMMMMMMM

- The first line gives a description of the protein and its content may be ignored.
- The second line gives the beginning of the primary structure of the protein, with letters encoding the amino acids. Note that the first amino acid is almost always M (methionine), as it corresponds to the start codon.
- The third line is an annotation corresponding to the localization of the cleavage site. For each amino acid on the previous line, a letter indicates whether this amino acid belongs to the signal peptide (annotated with a letter S) or to the mature protein (annotated with a letter M). The first amino acid of the mature protein is annotated with a letter C.

## References

- [1] Gunnar van Heine. A new method for predicting signal sequence cleavage sites. *Nucleic Acids Research* 14(11):4683–4690, 1986.
- [2] Jean-Philippe Vert. Support Vector Machines prediction of signal peptide cleavage site using a new class of kernels for strings. In: *Proceedings of the 7th Pacific Symposium on Biocomputing*, pp. 649–660, 2002.
- [3] [https://en.wikipedia.org/wiki/Additive\\_smoothing](https://en.wikipedia.org/wiki/Additive_smoothing)
- [4] [https://en.wikipedia.org/wiki/Substitution\\_matrix](https://en.wikipedia.org/wiki/Substitution_matrix)