

Evolving Training Sets for Peptide Discrimination via Evolutionary Algorithms

Loek Gerrits (s1032343), Evangelos Spithas (s1125593), Bart van Nimwegen

June 11, 2025

1 Introduction

In human bodies, the immune system is responsible for detecting pathogens and exterminating them. One of the tools at its disposal, is the T cells. T cells grow in the thymus, and each one is responsible for identifying and attacking specific cells. However, that means that it is possible for healthy cells to be attacked as well. In order to mitigate this, human cells are presented to the T cells, and the T cells that attack them are eliminated. This process is called Negative Selection (NS). In practice, the amount of possible cells is very large, and as a result only a sub-set of self peptides are presented in the T cells.

Artificial Immune Systems (AIS) are systems that draw inspiration from the human immune system, similarly to how Neural Networks (NNs) are inspired by the human nervous system. In this project, we would like to explore how we can find optimal subsets to train an NS algorithm for peptide selection. We will do this using 3 different methods, randomly sampling a subset, using a greedy algorithm and an evolutionary algorithm to produce them. Afterwards, we will train the negative selection algorithm with each one of them, and evaluate its performance against different sets of harmful peptides, such as HIV and ebola cells.

2 Methods

2.1 Datasets

2.1.1 Random

2.1.2 Greedy Algorithm

2.1.3 Evolutionary Algorithm

AA composition

AA frequency

Exchangeability

2.2 Negative Selection

In order to run the negative selection experiments we will use the jar that is provided by <https://johannes-textor.name/negativeselection.html>. The training of the algorithm will take place with the three training sets as described in section 2.1. All peptides have a fixed length of 6, and we will evaluate the effectiveness of the NS algorithm by experimenting with a variable length of contiguous selectors, ranging between 1 and 6.

In order to assess the performance of the NS, we will use as test sets, a subset of the self peptides, with size 1216, and a number of anomalous peptides that belong to the following categories: ebola, hepatitis b... For each peptide in the test datasets, the NS algorithm can return either the number of patterns in the repertoire that match it, or the normalised value $\log_2(1 + x)$, where x is the number of matching patterns. We will use this normalised value, as the number of selectors can be unwieldy to work with. Afterwards, we can classify each peptide as anomalous or self peptide, if its score exceeds the value r . After classifying the data, we will estimate the receiver operating characteristic curve (AUC) to quantify how well our estimated datasets can optimise the NS algorithm.

2.3 Experiment

2.4 Analysis

3 Results

4 Discussion

5 Conclusion

The end.