

# *Projekt lab notes*

2016

## *Project*

Exploration and implementation of algorithm for classification of signal-peptides.

---

*19 december 2016*

We have decided to mainly focus on implementing a classifier using some sort of HMM approach. If there is time we will also try to do something using an RNN. The first step will be writing code to handle and import the data.

*20 december 2016*

Today we implemented the first attempt at using an HMM. We used the hidden states from the data to train two models. One on the positive data and one on the negative data. We then used these to score the data points in the test set, choosing the model that had the highest probability score.

```
Number of positive samples labeled positive 528, total number of positive 528
Number of negative samples labeled negative 118, total number of negative 534
precision: 55.932203389830505
recall: 100.0
accuracy: 60.8286252354049
avrg positive neg prob: -4.320928254863987
avrg positive pos prob: -1.8246559833753646
avrg positive neg prob: -3.372393241005475
avrg positive pos prob: -3.1031834877775832
```

Figure 1: Results from first run.

This approach did not fare so well. It seems that the mean probability of the negative model is much lower, creating a classifier that is overly prone to positive classification.

We will now try a different approach, where we instead train a model on all the samples, and have it predict a hidden state sequence for the given protein sequence. We then use the hidden state sequence to predict the class of the data by looking for "C".

*21 december 2016*

Today we finished an single HMM approach to the problem. In this we use the hidden-state data to create a markow model for all the peptides in the training data. Then to classify new sequences we

use the model to predict the the hidden state of the sequence, and then look at the produced hidden-state to decide if the sequence is a signal peptide.

```

Evaluated on full data set:
Predicting.
Done.
-----
All true 127
All positive 128
True positive 117
True negative 128
Precision 0.9140625
Recal 0.92125984252
Accuracy on test data: 92.11%
-----
Evaluated on non-tm:
Predicting.
Done.
-----
All true 120
All positive 118
True positive 110
True negative 106
Precision 0.932203389831
Recal 0.916666666667
Accuracy on test data: 92.31%
-----
Evaluated on tm:
Predicting.
Done.
-----
All true 7
All positive 10
True positive 7
True negative 22
Precision 0.7
Recal 1.0
Accuracy on test data: 90.62%
-----

```

Figure 2: Stats from the second run

22 december 2016

We are now trying to use our model to analyze the proteome. We realized that our model had no way of handling errors in the data, or '\*' and had to adjust for this.

---