

Applied Machine Learning - Exercise 8 (15.06.2017)

Ben Wulf, Lie Hong, Amnon Bleich

Task 1)

We should use the 'clone mix' dataset. Therefore, we used reads which has the prefix 'acgagtgcgt'. We discharge any read which will not exactly map to this prefix. We count the occurrence of every sequence identical copy of the reads. Later, these counts are used to get the U in the EM.

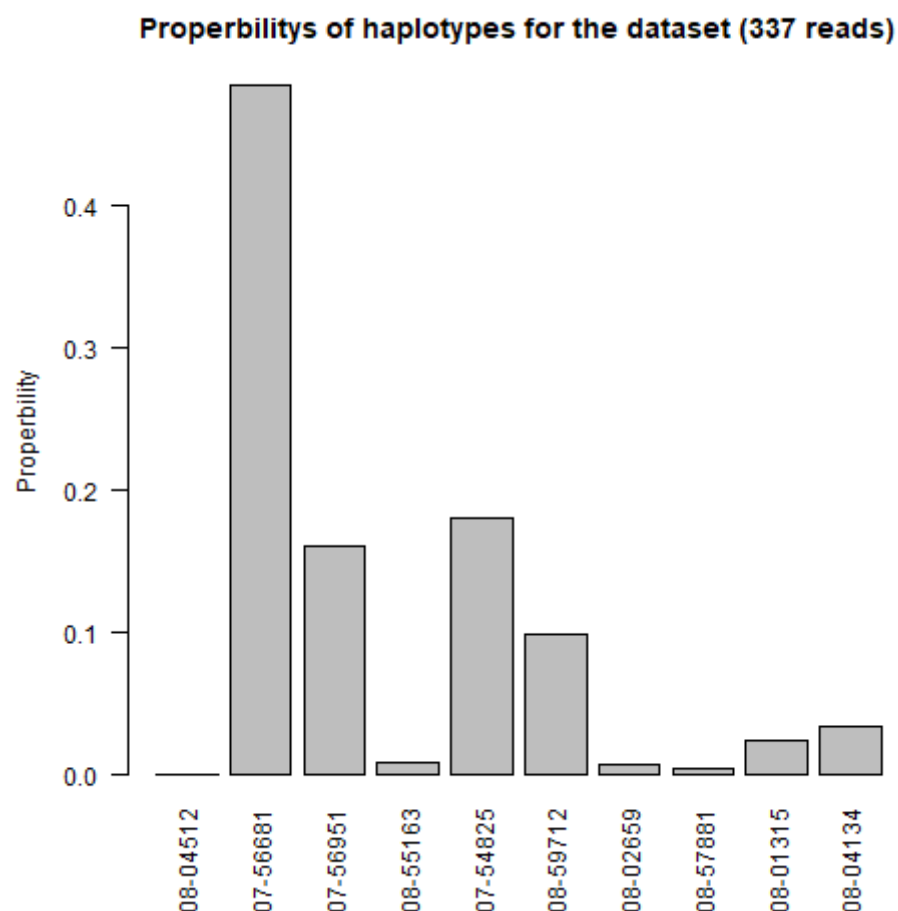
Task 2)

We create a boolean matrix, where every (unique sequence) read (r) is a row and every haplotype (h) is a column. An entry in this matrix is true if read R_r is in sequence identical content in the haplotype H_h, otherwise the entry is false. We create a subset of this matrix which will only contain sequences (reads) which will occur in at least one haplotype h. We also subset the counting vector of the read sequences to sequences that occur in at least one haplotype h. 337 different read sequences will remain.

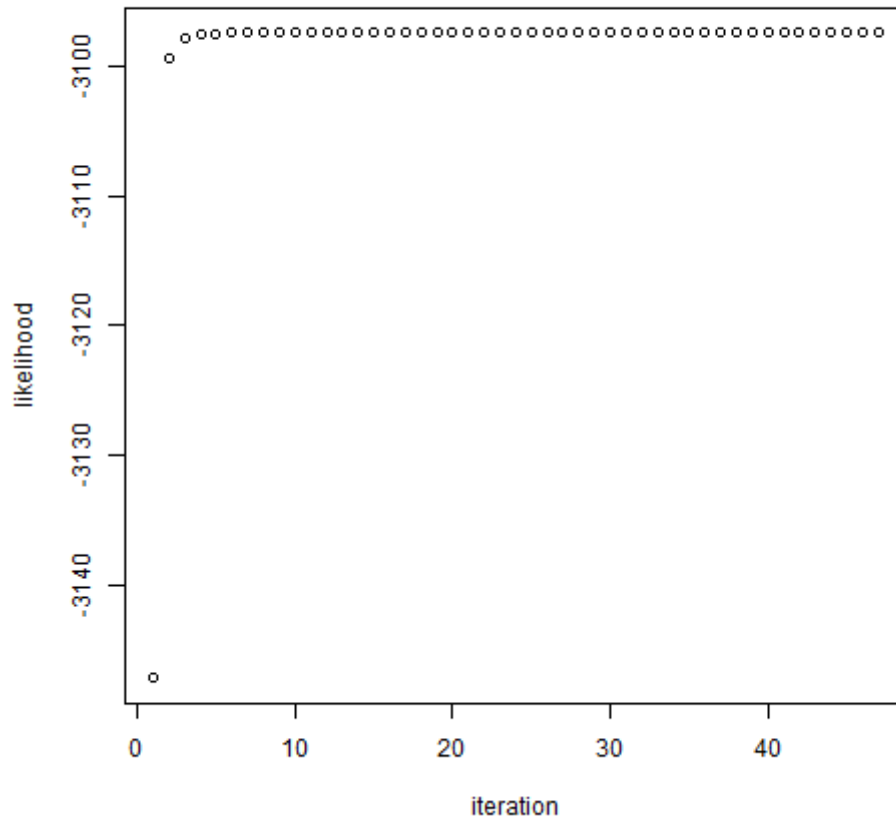
Task 3)

We were not sure what the U_r really means. We assume that U_r is the count of times a read sequence occurs in the dataset.

We decide to break the iterations for the EM if the changes of every p is smaller than a given epsilon. In the try with the full dataset we are using an epsilon of 10^{-10} . It seems that with a bigger epsilon the EM will never break.



We can see that haplotype 07-56681 is most likely for the given dataset and 08-04512 is extreme unlikely.



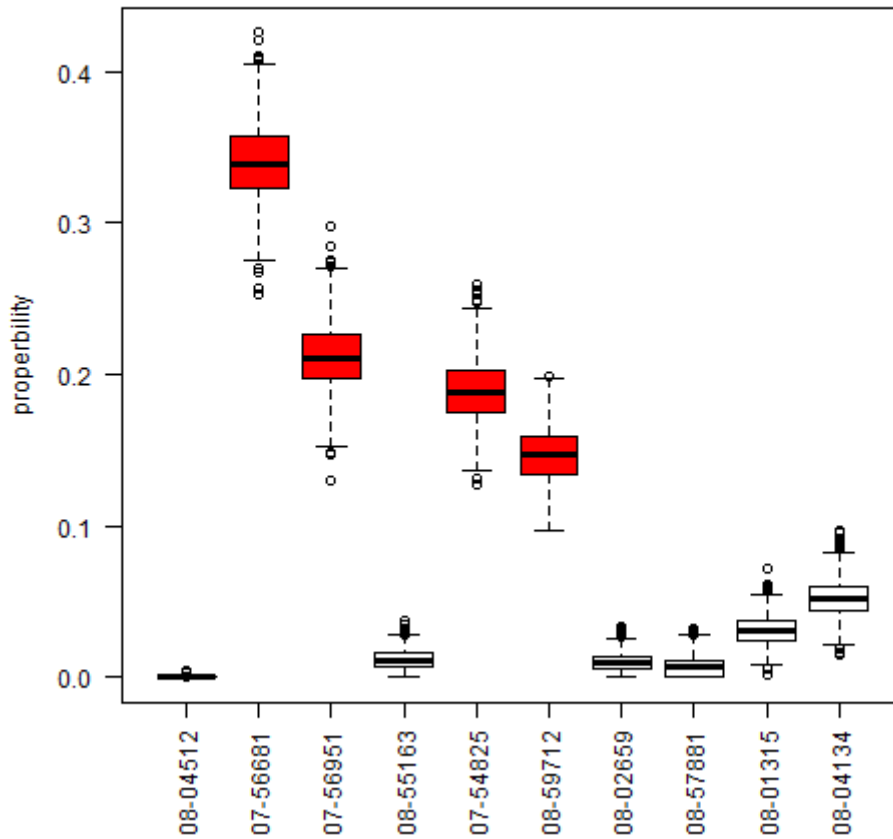
The likelihoods converging rapidly against the maximum and after the 5th iteration they are quite similar.

Task 4)

In each bootstrap iteration, we use again only unique read sequences for the matrix of read occurrences in the different haplotypes. The U vector is resampled for every iteration according to the sequences which are chose more than once. For time reasons, we decide to use a smaller epsilon of 10^{-5} in the bootstrap to break the EM iterations.

value			
haplotypes	mean	conf_intervall_low_95	conf_intervall_up_95
08-04512	8.816526e-06	-0.001406199	0.021233118
07-56681	-1.233802e-05	0.024039155	0.007229426
07-56951	9.037720e-09	0.188263546	-0.006528619
08-55163	3.406675e-01	0.149360962	0.016165533
07-54825	2.881505e-01	0.229167088	0.031193799
08-59712	3.869338e-01	0.146715937	0.011645883
08-02659	2.114684e-01	0.108908510	0.047215959
08-57881	1.676390e-01	0.183858822	0.052584845
08-01315	2.523766e-01	0.010170644	0.029184517
08-04134	1.169705e-02	-0.003134301	0.077082815

Task 5)



Which haplotypes are most dominant?

The most dominant haplotype is the second (07-56681), but also the 3rd, the 5th and the 6th haplotypes are common. The rest less common.

How reliable are results based on the bootstrapping?

The overlap between the boxes are not big. That means, that rank of the haplotypes driven by the bootstrap results should lead to a stable chance.

How common is a change in the order of importance of the haplotypes?

The change in the top 4 haplotypes is uncommon, because they differ allot in their distribution, but it is possible that the 5th haplotype is more likely than the 3rd. In the rest of the haplotypes are changes extremely common, because their distributions do not really differ.