

Homework 2

Björn Bebensee (2019–21343)
Machine Learning in Bioinformatics

June 10, 2020

1 Multiple Sequence Alignment using MUSCLE

First, in order to do further analysis with the given ten SARS-CoV-2 genome sequences and to generate the pHMM, we have to obtain a multiple sequence alignment. In order to align these sequences we can use MUSCLE [1]. Using MUSCLE we simply input the source sequences in .fasta format and run the command seen in Figure 1.

```
haean 福 ~/Documents/uni/Lecture Notes/2020 Spring/Machine Learning in Bioinformatics/HW2 > bf68553|master⚡
|9044 ± : muscle -in mlbio_input_seq.fasta -out input_seq.msa -clwstrict

MUSCLE v3.8.1551 by Robert C. Edgar

http://www.drive5.com/muscle
This software is donated to the public domain.
Please cite: Edgar, R.C. Nucleic Acids Res 32(5), 1792-97.

mlbio_input_seq 10 seqs, lengths min 29724, max 29903, avg 29813
00:00:00      3 MB(0%) Iter  1  100.00% K-mer dist pass 1
00:00:00      3 MB(0%) Iter  1  100.00% K-mer dist pass 2
00:02:58    1111 MB(6%) Iter  1  100.00% Align node
00:02:58    1111 MB(6%) Iter  1  100.00% Root alignment
00:04:56    1120 MB(7%) Iter  2  100.00% Refine tree
00:04:56    1120 MB(7%) Iter  2  100.00% Root alignment
00:04:56    1120 MB(7%) Iter  2  100.00% Root alignment
00:09:56    1120 MB(7%) Iter  3  100.00% Refine biparts
```

Figure 1: Command used to align the sequences

The output of MUSCLE is a multiple sequence alignment that can be seen in Figure 2, here in CLUSTAL W format. We can see that these sequences have now been aligned. Interestingly we can notice right away that the five sequences from Germany, namely MT358639, MT358640, MT358641, MT358642 and MT358643, have been perfectly aligned to each other at the beginning while the beginning of the five sequences from the USA, namely MT370833, MT370834, MT370835, MT370836 and MT370837, have been perfectly aligned to each other as well but offset by a little from the beginning of the sequences from Germany, meaning the American sequences do not contain this part at the beginning.

```

1 EUSTAL W (1.81) multiple sequence alignment
2
3 MT358641.1 ATTAAGGTTTACCTTCCAGGTAAACAAACCAACCAACTTCGATCTCTTAGATCT
4 _____ GATCT
5 MT370833.1 _____ GATCT
6 MT358639.1 ATTAAGGTTTACCTTCCAGGTAAACAAACCAACCAACTTCGATCTCTTAGATCT
7 MT358640.1 ATTAAGGTTTACCTTCCAGGTAAACAAACCAACCAACTTCGATCTCTTAGATCT
8 MT358642.1 ATTAAGGTTTACCTTCCAGGTAAACAAACCAACCAACTTCGATCTCTTAGATCT
9 MT370834.1 _____ GATCT
10 MT370836.1 _____ GATCT
11 MT358643.1 ATTAAGGTTTACCTTCCAGGTAAACAAACCAACCAACTTCGATCTCTTAGATCT
12 MT370835.1 _____ GATCT
13 *****
14
15 MT358641.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
16 MT370833.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
17 MT370837.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
18 MT358639.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
19 MT358640.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
20 MT358642.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
21 MT370834.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
22 MT370836.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
23 MT358643.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
24 MT370835.1 GTTCTAAACGAACCTTAAACCTGGCTGTCACTCGCGTCATCCTTAGTCAC
25 *****
26
27 MT358641.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
28 MT370833.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
29 MT370837.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
30 MT358639.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
31 MT358640.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
32 MT358642.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
33 MT370834.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
34 MT370836.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
35 MT358643.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
36 MT370835.1 CACCGAGTATAATTAAACTAAATTACTGGTTGACAGGACACGAGTAACTCGCTATC
37 *****
38
39 MT358641.1 TTCTGCAGGCTGTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
40 MT370833.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
41 MT370837.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
42 MT358639.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
43 MT358640.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
44 MT358642.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
45 MT370834.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
46 MT370836.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
47 MT358643.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
48 MT370835.1 TTCTGCAGGCTGCTTACGGGTTCTCGCGTGTGTCAGCGGATCATCACGACATCTAGGTT
49 *****
50
51 MT358641.1 CGTCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
52 MT370833.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
53 MT370837.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
54 MT358639.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
55 MT358640.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
56 MT358642.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
57 MT370834.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
58 MT370836.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
59 MT358643.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
60 MT370835.1 TGTCCGGGTGACCGAAAGTAAGATGGAGACCCCTGTCCTGGTTAACAGGAAAC
61 *****
62
63 MT358641.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
64 MT370833.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
65 MT370837.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
66 MT358639.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
67 MT358640.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
68 MT358642.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
69 MT370834.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
70 MT370836.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
71 MT358643.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
72 MT370835.1 ACACGTCAACTCAGTTGCCGTGTTTACAGGTTCGCGACGTGCTGTCAGTGCTTGG
73 *****
74
75 MT358641.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
76 MT370833.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
77 MT370837.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
78 MT358639.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
79 MT358640.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
80 MT358642.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
81 MT370834.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
82 MT370836.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
83 MT358643.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
84 MT370835.1 AGACTCGTGGAGGAGCTTACAGAGGCCGTCACACATTAAAAGATGGCACTTGG
85 *****
86
87 MT358641.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
88 MT370833.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
89 MT370837.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
90 MT358639.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
91 MT358640.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
92 MT358642.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
93 MT370834.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
94 MT370836.1 CTTAGTAGAAGTGAAGAAAGGGTTTGCCTCAACTTGAAACGCCCTATGTGTTCATCA
"input_seq.msa" 5989L, 422796C

```

Figure 2: Resulting multiple sequence alignment

2 Building a profile HMM using HMMER

Now that we have obtained a multiple sequence alignment for our sequences, we can run the `hmmbuild` command in HMMER [2] as seen in Figure 3. As we can see there are ten aligned sequences in the input multiple sequence alignment (`nseq`) with 29 903 aligned columns (`alen`). The resulting output profile by HMM has the same amount of consensus positions (`mlen`) and a relative entropy per position of 0.622 bits (`re/pos`).

```
haean 福 ~/Documents/uni/Lecture Notes/2020 Spring/Machine Learning in Bioinformatics/HW2 > bf68553|master⚡
|9058 ± : hmmbuild covid.hmm input_seq.msa
# hmmbuild :: profile HMM construction from multiple sequence alignments
# HMMER 3.3 (Nov 2019); http://hmmer.org/
# Copyright (C) 2019 Howard Hughes Medical Institute.
# Freely distributed under the BSD open source license.
# -----
# input alignment file:           input_seq.msa
# output HMM file:              covid.hmm
# -----
# idx name          nseq  alen  mlen    W eff_nseq re/pos description
#--- -----
1   input_seq      10 29903 29903 31848     1.70  0.622
# CPU time: 83.77u 0.20s 00:01:23.97 Elapsed: 00:01:24.34
```

Figure 3: Command to build the HMM

The resulting profile hidden Markov model can be seen in Figure 4. In this file we can see the metadata section and below it the actual HMM parameters. More specifically, these model parameters are negative log probabilities that describe the transition probabilities between the different states (that can be seen at the top of the parameter section). Each match state (enumerated on the left side from 1 to MAXL, here 29 903) contains three lines: match emission probabilities, insert emission probabilities and finally transition probabilities between the states.

```

1 HMMER3/f [3.3 | Nov 2019]
2 NAME input_seq
3 LENG 29903
4 MAXL 31848
5 ALPH DNA
6 RF no
7 MM no
8 CONS yes
9 CS no
10 MAP yes
11 DATE Mon Jun 8 00:16:03 2020
12 NSEQ 10
13 EFFN 1.704102
14 CKSUM 3798645073
15 STATS LOCAL MSV -18.4306 0.69322
16 STATS LOCAL VITERBI -18.9873 0.69322
17 STATS LOCAL FORWARD -11.0465 0.69322
18 HMM      A      C      G      T
19          m->m  m->i  m->d  i->m  i->i  d->m  d->d
20 COMPO  1.26776 1.56207 1.47607 1.27187
21      1.38629 1.38629 1.38629 1.38629
22      0.31396 3.66461 1.41124 1.46634 0.26236 0.00000   *
23      1 0.62298 2.01106 1.78035 1.82493   1 a - - -
24      1.38629 1.38629 1.38629 1.38629
25      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
26      2 1.85387 1.73195 1.88786 0.66352   2 t - - -
27      1.38629 1.38629 1.38629 1.38629
28      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
29      3 1.85387 1.73195 1.88786 0.66352   3 t - - -
30      1.38629 1.38629 1.38629 1.38629
31      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
32      4 0.62298 2.01106 1.78035 1.82493   4 a - - -
33      1.38629 1.38629 1.38629 1.38629
34      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
35      5 0.62298 2.01106 1.78035 1.82493   5 a - - -
36      1.38629 1.38629 1.38629 1.38629
37      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
38      6 0.62298 2.01106 1.78035 1.82493   6 a - - -
39      1.38629 1.38629 1.38629 1.38629
40      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
41      7 1.90507 2.16817 0.51598 1.96692   7 g - - -
42      1.38629 1.38629 1.38629 1.38629
43      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
44      8 1.90507 2.16817 0.51598 1.96692   8 g - - -
45      1.38629 1.38629 1.38629 1.38629
46      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
47      9 1.85387 1.73195 1.88786 0.66352   9 t - - -
48      1.38629 1.38629 1.38629 1.38629
49      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
50      10 1.85387 1.73195 1.88786 0.66352   10 t - - -
51      1.38629 1.38629 1.38629 1.38629
52      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
53      11 1.85387 1.73195 1.88786 0.66352   11 t - - -
54      1.38629 1.38629 1.38629 1.38629
55      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
56      12 0.62298 2.01106 1.78035 1.82493   12 a - - -
57      1.38629 1.38629 1.38629 1.38629
58      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
59      13 1.85387 1.73195 1.88786 0.66352   13 t - - -
60      1.38629 1.38629 1.38629 1.38629
61      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
62      14 0.62298 2.01106 1.78035 1.82493   14 a - - -
63      1.38629 1.38629 1.38629 1.38629
64      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
65      15 1.91393 0.66933 1.96289 1.60945   15 c - - -
66      1.38629 1.38629 1.38629 1.38629
67      0.06777 3.41842 3.41842 1.46634 0.26236 1.54840 0.23900
68      16 1.91393 0.66933 1.96289 1.60945   16 c - - -

```

"covid.hmm" 89732L, 5293424C

Figure 4: Excerpt from the resulting HMM

3 Querying the profile HMM

We can now use the obtained profile HMM to query it against sequence databases. In our case we have three such query databases. To run a search query of the profile HMM against the sequence database we use the `hmmsearch` command as can be seen in Figure 5. We redirect the shell output to another file to save it.

```
haean 福 ~/Documents/uni/Lecture Notes/2020 Spring/Machine Learning in Bioinformatics/HW2 > bf68553|master⚡
|9066 ± : hmmsearch covid.hmm query/mlbio_query_1_seq.fasta > query_1_seq.out
```

Figure 5: Command to query the HMM with a query database

We can repeat this for each of the query databases and obtain results that can be seen in Figure 6, Figure 7 and Figure 8. Each of these sequence databases only contains a single sequence so we only see a single sequence in the resulting ranking for each of the queries (it would have been possible to query for all of these at once, see Figure 9). The most important value to look at for us is the E-value, which represents the expected number of false-positives, and thus the lower the E-value ($\ll 1$) the more likely our profile sequences are homologous to the sequence in the query database. It is thus a measure of statistical significance. We can see for all three of our queries that the E-value is 0, meaning that these sequences are highly likely to be homologous (which is not very surprising since they all are SARS-CoV-2 sequences as well). We can also look at the bit score: the higher the bit score, the more significant the query database is to our profile. For all three queries we observe a bit score $> 40\,000$, again meaning the sequences are very significant. In Figure 9 we can easily see that the sequence MT459989 is most significant as it has the highest bit score but the difference in score is very small; all of the sequences are highly significant.

Lastly, we can take a look at the “best 1 domain” scores and E-values and see that they are nearly identical to the full sequence values, meaning that our profile sequence only consists of a single domain that was found in the query sequence (if the full sequence score was high but the best single domain score was low this could mean we are dealing with a repetitive sequence).

Figure 6: Result of query 1

Figure 7: Result of query 2

Figure 8: Result of query 3

```

1 # hmmsearch :: search profile(s) against a sequence database
2 # HMMER 3.3 (Nov 2019); http://hmmer.org/
3 # Copyright (C) 2019 Howard Hughes Medical Institute.
4 # Freely distributed under the BSD open source license.
5 # -----
6 # query HMM file: covid.hmm
7 # target sequence database: query/query_all.fasta
8 # -----
9
10 Query:      input_seq [M=29903]
11 Scores for complete sequences (score includes all domains):
12   --- full sequence ---   best 1 domain --- #dom-
13   E-value  score  bias   E-value  score  bias   exp N Sequence  Description
14   -----  -----  -----   -----  -----  -----   -----  -----  -----
15   0 40346.3 1020.5      0 40346.3 1020.3   1.0  1  MT459989.1 |Severe acute respiratory syndrome coronavirus 2
16   0 40319.0 1013.3      0 40318.8 1013.4   1.0  1  MT291831.1 |Severe acute respiratory syndrome coronavirus 2
17   0 40241.2 1009.7      0 40241.0 1009.7   1.0  1  MT077125.1 |Severe acute respiratory syndrome coronavirus 2
18

```

Figure 9: Results when querying all databases at once

References

- [1] Edgar, Robert C. "MUSCLE: a multiple sequence alignment method with reduced time and space complexity." *BMC Bioinformatics* 5.1 (2004): 113.
- [2] Eddy, Sean. "HMMER user's guide." *Department of Genetics, Washington University School of Medicine* 2.1 (1992): 13. See: