

Genome Sequence Analysis: Assignment 1

University of Cambridge

Henrik Åhl

January 27, 2017

Preface

This is an assignment report in connection to the *Genome Sequence Analysis* module in the Computational Biology course at the University of Cambridge, Michaelmas term 2016. All related code is as of January 27, 2017 available per request by contacting hpa22@cam.ac.uk. The approach is done after the guidelines specified by Rabiner [1], and as summarized by Shen [2]. All code is implemented using Python, and in particular by using the `numpy` package.

Exercises

1

```
1 # Get index for letter (symbol) occurring.
def get_ind(probs):
3     rand = random.random()
    for ii in xrange(len(probs)):
5         if rand < probs[ii]:
            return ii
7
9 # Read in transition matrix and initial distribution from files.
# If sumrows == TRUE, sum the rows
10 def read_and_output(mat_file, idist_file, seq_length, sumrows, transpose):
    # Read in data
13     trans_mat = np.loadtxt(open(mat_file, "r"))
    init_dist = np.loadtxt(open(idist_file, "r"))
15
    # Act on the optional things
17     if transpose:
        trans_mat = np.transpose(trans_mat)
19     if sumrows:
        trans_mat = np.cumsum(trans_mat, 1)
21
    # Fill out the output probabilistically
23     output = seq_length * [None]
    output[0] = np.random.choice(np.arange(0, len(trans_mat)), p=list(init_dist))
25     for ii in xrange(1, seq_length):
        output[ii] = get_ind(trans_mat[output[ii - 1], :])
27
    "".join(str(s) for s in output)
29     return output
31
mat_file = "/home/henrik/compbio/src/assignments/gsa1/report/mat_file.dat"
33 idist_file = "/home/henrik/compbio/src/assignments/gsa1/report/idist_file.dat"
output = read_and_output(mat_file, idist_file, 115, True, False)
```

Evoking the `print` command we get the following:

[illegible]

Our program is thus able to produce an output corresponding to that of a Markov chain.

2

For this exercise, as we are only to read a single Markov chain, the initial distribution will correspond to a 100 % probability of ending up in the first state in the sequence.

```

3 # Get transition matrix and initial distribution from one (!) sequence
4 def infer_mat_and_init(sequence):
5     # Get some data and create the matrices we need.
6     elements = np.unique(list(sequence))
7     trans_mat = np.zeros((len(elements), len(elements)))
8     init_dist = np.zeros(len(elements)) # This is specific to this assignment
9     init_dist[int(sequence[0])] = 1 # Count the first one
10
11     # Count the occurrences of all events
12     for ii in xrange(1, len(sequence)):
13         trans_mat[int(sequence[ii - 1]), int(sequence[ii])] += 1
14
15     # What are the probabilities?
16     for ss in xrange(len(trans_mat)):
17         trans_mat[ss,] /= np.sum(trans_mat[ss,])
18
19     return trans_mat, init_dist
20
21 trans_mat, init_dist = infer_mat_and_init("
01010101011110100101011001010010101010101010010101")

```

Again, looking at the inferred transfer matrix and initial distribution tells us that our program appears to give an output corresponding to our expectations:

$$A = \begin{pmatrix} 0.15 & 0.85 \\ 0.85 & 0.15 \end{pmatrix}, \quad \mu^0 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

3

```

3 # Parameters
4 S = np.array([0, 1])
5 V = np.array([0, 1, 2, 3, 4])
6 A = np.matrix("0.8, 0.2;"
7               "0.1, 0.9")
8 B = np.matrix("0.2, 0.5, 0.2, 0.1, 0.0;"
9               "0.0, 0.1, 0.4, 0.4, 0.1")
10 mu0 = np.array([.5, .5])
11
12 # Emit a sequence given some transition and emission matrices, as well as
13 # an initial distribution.
14 def emit_sequence(seq_length, trans_matrix, emiss_matrix, init_dist):
15     c_trans_mat = np.cumsum(trans_matrix, 1)
16     c_emiss_mat = np.cumsum(emiss_matrix, 1)
17
18     emit = np.zeros(seq_length)
19     hidden = np.zeros(seq_length)
20     hidden[0] = np.random.choice(np.arange(0, len(trans_matrix)), p=list(init_dist))
21     emit[0] = get_ind((c_emiss_mat[int(hidden[0]), :]).tolist()[0])

```

```

22  # Parse sequence
    for ii in xrange(1, seq_length):
24      hidden[ii] = get_ind((c_trans_mat[int(hidden[ii - 1]), :]).tolist()[0])
        emit[ii] = get_ind((c_emiss_mat[int(hidden[ii]), :]).tolist()[0])
26
    emit = "".join(str(int(x)) for x in emit)
28    hidden = "".join(str(int(x)) for x in hidden)
    return emit, hidden
30
32 emit, hidden = emit_sequence(115, A, B, mu0)

```

The result of our emitted sequence, as well as the underlying state, can be seen in fig. 3, and appears to correspond well to the expected behaviour.

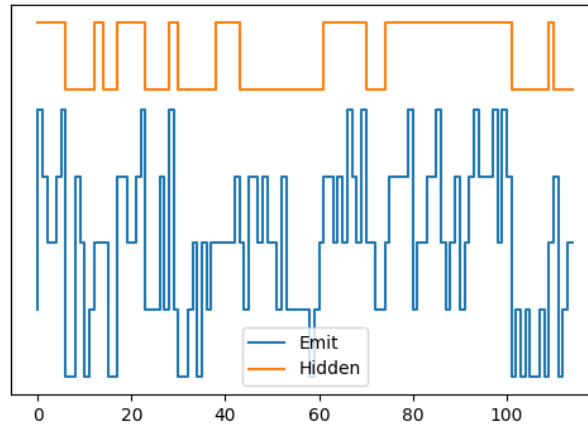


Figure 1: Output from the program emitting a sequence according to the designated parameters. 0-valued states tend to produce the lowest emission values on average, whereas the 1-valued states behave in the opposite manner.

4

We define a method for reading in a sequence from a file (which has been produced with our previously defined emission function).

```

# Read in sequence from file.
2 def read_sequence(file, skipFirst):
    seq = ""
4     with open(file, 'r') as in_file:
        if skipFirst:
6             in_file.readline()
            seq = in_file.read().replace('\n', '')
8     return seq

10
# Use the forward algorithm to determine (log) likelihood of sequence
12 def forward(sequence, trans_matrix, emiss_matrix, hidden_states, init_dist):
    # Function for getting normalisation constant
14     def get_constant(a):
        sum = np.sum(a)
16         return 1.0 / len(trans_matrix) if sum == 0 else 1.0 / sum

18     # Pre-allocate
    fw = np.zeros((len(hidden_states), len(sequence)))
20     constants = []

```

```

22     # Initialise
    for ss in xrange(len(hidden_states)):
24         fw[ss, 0] = init_dist[ss] * emiss_matrix[ss, int(sequence[0])]
    const = get_constant(fw[0, :])
26     constants.append(const)

28     # Scale accordingly
    for ss in xrange(len(hidden_states)):
30         fw[ss, 0] *= const

32     # Now go through the rest of the sequence
    probs = np.zeros(len(hidden_states))
34     for ii in xrange(1, len(sequence)):
        for ss in xrange(len(hidden_states)):
36             probs[ss] = np.sum(
                (fw[substate, ii - 1] * trans_matrix[substate, ss] * emiss_matrix[ss, int(
38                 sequence[ii])]) for
                substate in xrange(len(hidden_states)))

40     # Scale again
    const = get_constant(probs)
42     constants.append(const)
    for ss in xrange(len(hidden_states)):
44         fw[ss, ii] = const * probs[ss]

46     # Sum up the probabilities; disregard if 0
    ln_prob = -np.sum([math.log(const) if const != 0 else 0 for const in constants])
48     return ln_prob, constants, fw

50 out_sequence, -- = emit_sequence(115, A, B, mu0)
seq_file = '/home/henrik/compbio/src/assignments/gsa1/random_output_sequence.dat'
52 with open(seq_file, "w") as f_out:
    f_out.write(out_sequence)
54     f_out.close()
that_same_sequence = read_sequence(seq_file, False)
56 prob, consts, fw = forward(that_same_sequence, A, B, S, mu0)

```

5

We use the *Saccharomyces cerevisiae* chromosome III genome from Ensembl [3]. The genome is binned by taking intervals of size 100, and subsequently identifying how many G or C bases are contained. In our implementation, bins are sequential and do not overlap.

Trying to approximate the correct distribution of categories, we label bins after percentage-wise GC content. Doing this rough fitting, the best fit appears to be found around $0 - 28.5\%$, $28.5 - 34.2\%$, $34.2 - 40.0\%$, $40.0 - 50.0\%$ and a $50 - 100\%$ split for label 1-5 respectively, with upper boundaries exclusive. As our genome is not evenly divisible by the bin size, the last bin will be misrepresentative. However, as our genome is large, and we thus will also have a large set of bins, the effects of this will be marginal. Thence we choose to disregard this effect.

```

# We cheat a little and put the remnants in a bin of their own, i.e. if the last bin only has
# e.g. 20 elements, the GC
2 # content in this will be "the number of GC bases in those 20 bases"/bin_size, which is wrong
# , but we'll have to
# live with it.
4 def calc_gc(sequence, bin_size):
    gc_cont = []
6     bin_size = float(bin_size)
    for ii in xrange(int(math.ceil(float(len(sequence)) / bin_size))):
8         ceil = int((ii + 1) * bin_size)
        if ceil > len(sequence):
10             ceil = len(sequence)
        bin = sequence[int(ii * bin_size): ceil]

```

```

12     G = bin.count("G")
13     C = bin.count("C")
14     gc_cont.append((G + C) / bin_size)
15     return gc_cont
16
17
18 # Read genome
19 sc_file = '/home/henrik/compbio/src/assignments/gsal/sc-gen.fa'
20 sc_gen = read_sequence(sc_file, True)
21 gc = calc_gc(sc_gen, 100)
22
23
24 # Relabel according to predefined cuts
25 def relabel(seq, cuts):
26     relab = []
27     for ii in seq:
28         for cut in xrange(len(cuts)):
29             if ii <= cuts[cut]:
30                 relab.append(cut)
31                 break
32
33     relab = "".join(str(x) for x in relab)
34     return relab
35
36 # Compare the two strings
37 relab = relabel(gc, [.285, .342, .40, .5, 1.]) # Split by percentages
38 model3_seq, model3_hidden = emit_sequence(len(relab), A, B, mu0)
39 log_p, const, fw_seq = forward(relab, A, B, S, mu0)

```

The final classification can be seen in fig. 2, where the two distributions are shown.

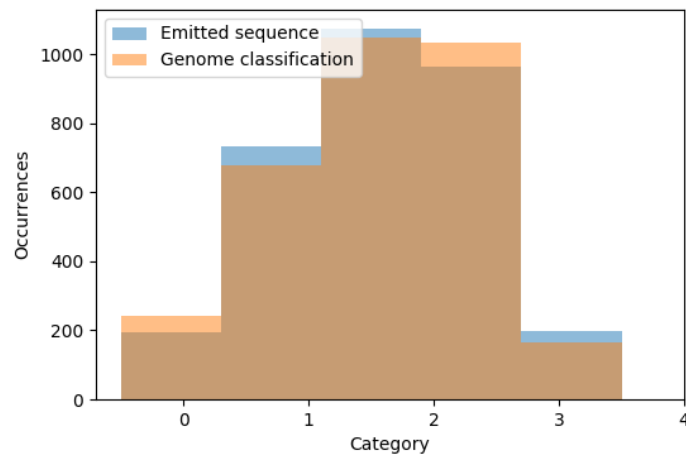


Figure 2: Distributions of emission values according to our given model and of the *Saccharomyces cerevisiae* genome after percentage-wise labelling

6

```

1 # We only use this for Baum-Welch, so we use the same constants as in our forward case
2 def backward(sequence, trans_matrix, emiss_matrix, hidden_states, init_dist, constants):
3     # Preallocate for our backwards sequence
4     bw = np.zeros((len(hidden_states), len(sequence)))
5
6     # Initialise

```

```

7     for ss in xrange(len(hidden_states)):
8         bw[ss, len(sequence) - 1] = 1.
9         bw[ss, len(sequence) - 1] = constants[len(sequence) - 1]

11    # Go through the rest of the sequence, continue from the back end
12    for ii in xrange(len(sequence) - 2, -1, -1):
13        b = np.zeros(len(hidden_states))
14        for ss in xrange(len(hidden_states)):
15            b[ss] = np.sum([(bw[substate, ii + 1] * trans_matrix[ss, substate] *
16                             emiss_matrix[substate, int(sequence[ii + 1])]) for substate in
17                                xrange(len(hidden_states))])

18    # Again, scale
19    for ss in xrange(len(hidden_states)):
20        bw[ss, ii] = constants[ii] * b[ss]

21    ln_prob = -np.sum([math.log(c) if c != 0 else 0 for c in constants])
22    return ln_prob, bw

25    # Now we can go onto estimating our matrices
26    def baum_welch(sequence, trans_matrix, emiss_matrix, hidden_states, init_dist):
27        # Get those lists we need
28        p1, constants, fw = forward(sequence, trans_matrix, emiss_matrix, hidden_states,
29                                     init_dist)
30        p1, bw = backward(sequence, trans_matrix, emiss_matrix, hidden_states, init_dist,
31                           constants)

32        # Initial distribution
33        for ss in xrange(len(hidden_states)):
34            init_dist[ss] = fw[ss, 0] * constants[0] * bw[ss, 0]
35        init_dist /= np.sum(init_dist)

36        # New state matrix
37        for ss in xrange(len(hidden_states)):
38            for substate in xrange(len(hidden_states)):
39                # Sum up the estimate, normalise and assign
40                new_est = np.sum([(fw[ss, ii] * bw[substate, ii + 1] * trans_matrix[ss, substate]
41                                     * emiss_matrix[
42                                         substate, int(sequence[ii + 1])]) for ii in xrange(len(sequence) - 1)])
43                norm = np.sum([(fw[ss, ii] * bw[ss, ii] / constants[ii]) for ii in xrange(len(
44                    sequence) - 1)])
45                trans_matrix[ss, substate] = new_est / norm if norm != 0 else 1. / len(
46                    hidden_states)

47        # New emission matrix (structure same as above)
48        for ss in xrange(len(hidden_states)):
49            for jj in xrange(emiss_matrix.shape[1]):
50                new_est = np.sum([(fw[ss, ii] * bw[ss, ii] / constants[ii] if int(sequence[ii])
51                                     == jj else 0) for ii in
52                                     xrange(len(sequence))])
53                norm = np.sum([(fw[ss, ii] * bw[ss, ii] / constants[ii]) for ii in xrange(len(
54                    sequence))])
55                emiss_matrix[ss, jj] = new_est / norm if norm != 0 else 1. / emiss_matrix.shape
56                [1]

57        return p1

58    change = 9999999999999999
59    prev = change
60    while change > 0.000001:
61        prob = baum_welch(relab, A, B, S, mu0)
62        change = prev - prob
63        prev = prob

```

We here achieve a steady increase in the log likelihood for the given sequence, showing that our model adapts

to the data given. Ultimately, our matrices take the shapes of

$$A = \begin{pmatrix} 0.91 & 0.09 \\ 0.05 & 0.95 \end{pmatrix}, \quad B = \begin{pmatrix} 0.01 & 0.05 & 0.21 & 0.59 & 0.14 \\ 0.11 & 0.31 & 0.40 & 0.18 & 0.00 \end{pmatrix}.$$

Our final log likelihood for the same situation evaluates as:

-4251.17034431

7

```

1  # Retrieve the most likely input hidden states given and output sequence and
  # the corresponding matrices.
3  def viterbi(sequence, trans_matrix, emiss_matrix, hidden_states, init_dist):
    # Define some things we're gonna use.
5     constants = np.zeros(len(sequence))
    solution = np.zeros(len(sequence))
7     dyn_mat = np.zeros((len(hidden_states), len(sequence)))
    store_path = np.zeros((len(hidden_states), len(sequence)))
9
    # Initialise rows
11    dyn_mat[:, 0] = init_dist * emiss_matrix[:, int(sequence[0])]
    constants[0] = 1.0 / np.sum(dyn_mat[:, 0])
13    dyn_mat[:, 0] *= constants[0]

15    # Fill out the table
    for ii in xrange(1, len(sequence)):
17        for ss in xrange(len(hidden_states)):
            probabilities = np.diag(dyn_mat[:, ii - 1]) * trans_matrix[:, ss] # Probability
                of coming from either state
19            store_path[ss, ii], dyn_mat[ss, ii] = max(enumerate(probabilities), key=
                itemgetter(1))
            dyn_mat[ss, ii] *= emiss_matrix[ss, int(sequence[ii])]
21
        # Now we scale
23        constants[ii] = 1.0 / np.sum([dyn_mat[ss, ii] for ss in xrange(len(hidden_states))])
        dyn_mat[:, ii] *= constants[ii]
25
    # Backtrack
27    solution[len(sequence) - 1] = dyn_mat[:, len(sequence) - 1].argmax() # last state
    for ii in xrange(len(sequence) - 1, 0, -1):
29        solution[ii - 1] = store_path[int(solution[ii]), ii]

31    return solution

33
35 model7_seq = sc_gen
gc = calc_gc(sc_gen, 100)

37 intab = "ATCG"
outtab = "0123"
39 trantab = maketrans(intab, outtab)
model7_seq = model7_seq.translate(trantab)
41
43 change = 9999999999999999
prev = change
while change > 0.000001:
45     prob = baum_welch(relab, A, B, S, mu0)
    change = prev - prob
47     prev = prob

49 solution = viterbi(relab, A, B, S, mu0)

```

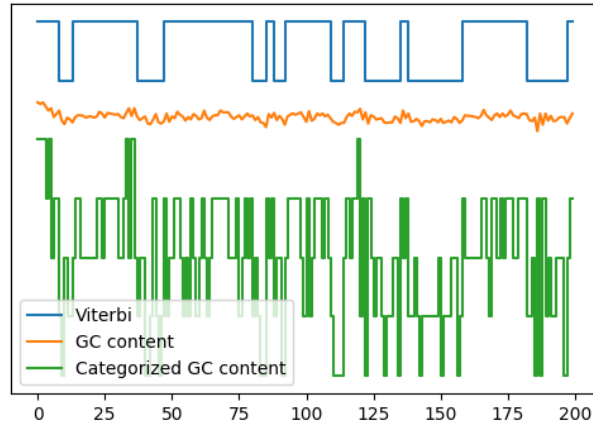


Figure 3: Inferred underlying hidden state sequence using the Viterbi algorithm, along with the corresponding categorized and non-categorized GC content under the previously determined binning scheme. Notably, the 1-valued states correspond well with the highest gc content, as expected from our matrices. Similarly, the lower state more consistently is affiliated with lower emission values.

8

Categorise according to a gaussian distribution with different means?