

# **Protein Sequence Reconstruction (Option 4: HMM Applications)**

Leon Zhang, Qiubai Yu

November 2019

# 1 Title of Project

Protein Sequence Reconstruction

## 2 Names and Roles of Each Teammate

### 2.1 Leon Zhang (leonz, 1664187)

1. Prepare protein sequence alignment data
2. Calculated the emission probability from the data
3. Implemented the Viterbi algorithm for the use of protein sequence reconstruction

### 2.2 Qiubai Yu (qiubay, 1663777)

1. Calculated the transition probability from the data
2. Implemented the Forward algorithm
3. Testing Code

### 2.3 Common contribution

1. Project idea and research
2. Final report

## 3 What the Program is Supposed to do

Proteins are macro molecules made my DNA, consisting of one or more long chains of amino acid residues. There are only 20 amino acids, but this limited number of amino acids are combined in infinite number of ways that enable biological functions such as catalysing metabolic reactions, transporting molecules, responding to stimuli, etc<sup>[2]</sup> to take place. In many cases, when DNA become mutated, they fail to produce expected protein sequences, which could potentially cause functional deficiency or abnormality in the organism. Many of today's most intractable diseases, including cancers, Alzheimer's disease, color blindness and so on, are caused by dramatically incorrect protein sequences<sup>[3]</sup>.

Protein sequence study is hot field. Protein sequence alignment is a common topic in bioinformatics. Sequence alignment tools are developed to find the differences between two proteins by guessing whether there is an amino acid deletion, insertion or substitution from one sequence to another. By comparing a functional and a nonfunctional mutated protein sequence, scientists are able to discover the missing or mutated amino acids, and even predict what should

the functional sequence should be.

Our goal is to use Hidden Markov Model to achieve protein sequence reconstruction (from a mutated sequence). Our program will take a mutated protein sequence (has deleted, mutated and inserted amino acids) and predict what the functional protein sequence is. More detailed approach and analysis is included in below sections.

## 4 Technique Used

We began by first downloading data from National Center for Biotechnology Information (NCBI). Our data are functional and nonfunctional TRPA1 protein sequences of mouse family. TRPA1 is a type of sensor protein that gives rise to sensory modalities, such as pain, cold and itch, and other protective responses. We chose to use TPRA1 protein because it is found in variety range of animals and its reasonable length (around 1000 amino acids). We collected these protein sequences from mouse family animals including: House mouse, Thirteen-lined ground squirrel, Chinese Hamster (*Cricetulus griseus*), Guinea pig (*Cavia porcellus*), Beaver (*Castor canadensis*), Hedgehog (*Echinops telfairi*), Ferret (*Mustela putorius furo*), and Chinchilla (*Chinchilla lanigera*).

We then used EMBOSS Needle<sup>[3]</sup>, a online protein alignment tool to compare the two sequences to figure out where are the deletions or insertions of amino acids. A sample image is as following:

After comparing, the EMBOSS Needle shows us there is 90.7 percent similarity between the two sequence (high because they are both from mouse family), and there is 34 deletions. Below the score section shows the actual differences of the two protein sequences. '-' represents a match, ':' represents a substitution, and '.' represents a deletion.

We then use the data to find out our transition and emission probabilities of our HMM Models which is shown below.

The transition probability is the probability of states within a sequence (Within a protein sequence each amino acid has a probability of going to another amino acid or a deletion state). The emission probability is the probability of states going reconstructed protein sequence.

The Viterbi algorithm we implemented uses these transition and emission probabilities to find the most probable functional protein sequence from the mutated sequence. The Viterbi algorithm works by continuously finding the most likely path up to each state using previously found best path. This algorithm is very efficient (runs in  $O(T \cdot N^2)$  where  $T$  is the length of observations, and  $N$  is the number of hidden states), since it reduces the computations all calculating all the path.

```
#####
# Program: needle
# Runday: Mon 2 Dec 2019 06:41:10
# Commandline: needle
#
# -auto
# -stdout
# -asequence emboss_needle-I20191202-064107-0840-76085776-p2m.asequence
# -bsequence emboss_needle-I20191202-064107-0840-76085776-p2m.bsequence
# -datafile EBLOSUM62
# -gapopen 10.0
# -gapextend 0.5
# -endopen 10.0
# -endextend 0.5
# -aformat3 pair
# -sprotein1
# -sprotein2
# Align_format: pair
# Report_file: stdout
#####

#-----
# Aligned_sequences: 2
# 1: EMBOS_001
# 2: EMBOS_001
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 1125
# Identity: 944/1125 (83.9%)
# Similarity: 1020/1125 (90.7%)
# Gaps: 34/1125 ( 3.0%)
# Score: 4979.5
#
#-----

EMBOS_001      1 mkrgrlrrillpserkevqgyvrygvgedmdeskesfkvdiedgmcrldef 50
EMBOS_001      1 mkrslrkmlrpge-kepqdvvyqgvddmdesksdfkvvfegnahrlqnf 49
EMBOS_001     51 iknrrklslkyedenlcplhhaaaegqvemeliingscsevlinmdgygn 100
EMBOS_001     50 iknrrklslkyddtnasplhhaaaegqvemktilsgsscavlnvmdygn 99
EMBOS_001    101 tplhcaaaekngvesvkflslsgganpnlrnrnmmsplhiavhgmynveikv 150
EMBOS_001    100 tplhwateknqvesvkflslsgganpnlrnrnmmsplhiiaaagmyneikv 149
EMBOS_001    151 ltehkatninlegengntalmstcakdnsealqillekgaklckskankwgd 200
EMBOS_001    150 ltehrtstinlegengntallttcakdnsealqillnkaklckskankwgd 199
EMBOS_001    201 ypvhqaafsgakkcmelilaygekngysrethinfvnhkkaasplhlavgs 250
EMBOS_001    200 fpvhqaafsgakkcmelilkhgeehgysrgshinfvnnkksplhlavgs 249
```

Figure 1: Differences between TRPA1 Sequence of House Mouse and Thirteen-lined ground Squirrel

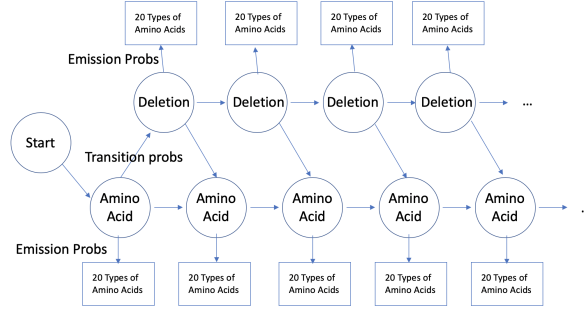


Figure 2: Hidden Markov Model for Protein Sequence Reconstruction

The Forward Algorithm is used to compute the joint probability of hidden state  $h_t$  and observations  $o_{1:t}$ . Before implementing the forward algorithm, being able to markov chain and conditional independent is needed<sup>[5]</sup>.

$$P(h_t, o_{1:t}) = \sum_{h_{t-1}} P(h_t, h_{t-1}, o_{1:t}), t \in [2, \text{inf}]$$

$$\begin{aligned}
&= \sum_{h_{t-1}} P(o_t|h_t, h_{t-1}, o_{1:t-1}) \cdot P(h_t|h_{t-1}, o_{1:t-1}) \cdot P(h_{t-1}, o_{1:t-1}) \\
&= \sum_{h_{t-1}} P(o_t|h_t) \cdot P(h_t|h_{t-1}) \cdot P(h_{t-1}, o_{1:t-1})
\end{aligned}$$

emission probability \* transition probability \* previously computed probability

For  $t = 1$ ,

$$P(h_1, o_1) = P(o_1|h_1) \cdot P(h_1)$$

emission probability \* initial probability

When implementing forward algorithm, dynamic programming is required. The forward algorithm stores previously computed probabilities and access them directly when they are needed.

## 4.1 Code

### 4.1.1 Generate Probability Matrices

```

1 # others observation
2 # mouse hidden states
3 # use others to compute transition probability
4 # compare and compute emission probability
5 init_dic = {}
6 for i in range(len(mouse_new)):
7     char = mouse_new[i]
8     if char in init_dic:
9         init_dic[char] += 1
10    else:
11        init_dic[char] = 1
12 #init_dic
13 summ = 0
14 for i in init_dic:
15     summ += init_dic[i]
16
17 for i in init_dic:
18     init_dic[i] = init_dic[i] / summ
19
20 init_dic
21
22 dic_1 = {}
23 for i in range(len(others_new) - 1):
24     first = others_new[i]
25     second = others_new[i + 1]
26     if first in dic_1:
27         dic_2 = dic_1.get(first)
28         if second in dic_2:
29             dic_2[second] += 1
30         else:
31             dic_2[second] = 1

```

```

32     else:
33         dic_1[first] = {}
34         dic_1[first][second] = 1
35
36 transition_key = ['-', 'a', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'k',
37 'l', 'm', 'n', 'p', 'q', 'r', 's', 't', 'v', 'w', 'y']
38 for key1 in dic_1:
39     dictionary = dic_1[key1]
40     for key2 in transition_key:
41         if not (key2 in dictionary):
42             dictionary[key2] = 0
43
44 # convert to probabilities
45 for key1 in dic_1:
46     dictionary = dic_1[key1]
47     summ = 0
48     for key2 in dictionary:
49         summ += dictionary[key2]
50     for key3 in dictionary:
51         dictionary[key3] = dictionary[key3] / summ
52 transition_dict = dic_1
53 transition_dict
54
55 emission_dict = {}
56 emission_keys = ['-', 'a', 'c', 'd', 'e', 'f', 'g', 'h', 'i', 'k',
57 'l', 'm', 'n', 'p', 'q', 'r', 's', 't', 'v', 'w', 'y']
58 for key in emission_keys:
59     emission_dict[key] = {}
60     for key2 in emission_keys:
61         emission_dict[key][key2] = 0
62
63 for i in range(len(mouse_new)):
64     mouse_aa = mouse_new[i]
65     others_aa = others_new[i]
66     if mouse_aa in emission_dict:
67         dic_2 = emission_dict.get(mouse_aa)
68         if others_aa in dic_2:
69             dic_2[others_aa] += 1
70         else:
71             dic_2[others_aa] = 1
72     else:
73         emission_dict[mouse_aa] = {}
74         emission_dict[mouse_aa][others_aa] = 1
75
76 for key in emission_dict:
77     curr_dict = emission_dict[key]
78     curr_total = 0
79     for key2 in emission_dict[key]:
80         curr_total += emission_dict[key][key2]
81     for key2 in emission_dict[key]:
82         emission_dict[key][key2] /= curr_total
83 emission_dict

```

#### 4.1.2 Forward Algorithm

```

1 # Forward Algorithm
2 # Compute the Pr(hidden.state:t, all_observations.1.to:t)
3 # denote hidden state h, observation o
4 # Pr(h:t, o.1:t) = sum of Pr(h:t, h:t-1, o.1:t)
5 # = Pr(o:t | h:t) * Pr(h:t | h:t-1) * Pr(h:t-1, o.1:(t-1)) for t ...
   in [2, inf)
6 # = emission_prob * transition_prob * previously_computed_prob
7 # for t = 1
8 # Pr(h.1, o.1) = Pr(o.1 | h.1) * Pr(h.1)
9 #               = emission_prob * initial_prob
10
11 import numpy as np
12 # obs {Array} - the observation sequence
13 # ep {nested dictionary} - emission probabilities
14 # tp {nested dictionary} - transition probabilities
15 # ip {dictioanry} - initial probabilties
16 # be able to return the joint probability of o.1:t and h:t
17 def forward(obs, ep, tp, ip):
18     num_obs = len(obs)
19     states = [s for s in emission_prob.keys()]
20     num_states = len(states)
21
22
23     # use a numpy matrix as a 2-D array to store probabilities
24     mat = np.zeros((num_states, num_obs))
25
26     # initialize the first column
27     first_ob = obs[0]
28     for i in range(num_states):
29         mat[i, 0] = ip[first_ob] * ep[states[i]][first_ob]
30
31     # compute the following columns recursively
32     for i in range(1, num_obs):
33         cur_obs = obs[i]
34         for j in range(num_states): # current hidden state s_j
35             summ = 0.0
36             for k in range(num_states): # previous hidden state s_k
37                 #print("i: ", i, "j: ", j, "k: ", k)
38                 summ += mat[k, i-1] * ep[states[j]][cur_obs] * ...
                    tp[states[k]][states[j]]
39             mat[j, i] = summ
40
41     # give out the predicted sequence
42     predicted_seq = []
43     for i in range(len(obs)):
44         cur_max = -1
45         index = -1
46         for j in range(len(states)):
47             if mat[j, i] > cur_max:
48                 cur_max = mat[j, i]
49                 index = j
50     predicted_seq.append(states[index])
51
52     return mat, predicted_seq

```

### 4.1.3 Viterbi Algorithm

```
1 def viterbi(transition_matrix, emission_matrix, initial_prob, ...
2     observations):
3     # Initialize returning objects
4     prediction_seq = []
5     prediction_prob = 0
6     # Initialize a dictionary for scores of each states
7     scores = {i : [] for i in list(transition_matrix.keys())}
8     # Fill in first scores column (start state to transition ...
9     states probability)
10    max_score = 0
11    best_key = ''
12    for key in scores:
13        if key in emission_matrix or key in initial_prob:
14            curr_observation = observations[0]
15            score = ...
16            emission_matrix[key][curr_observation]*initial_prob[key]
17            scores[key].append(score)
18            if score > max_score:
19                max_score = score
20                best_key = key
21            else:
22                scores[key].append(0.0)
23    prediction_seq.append(best_key)
24    # Update the rest of the scores
25    for i in range(1, len(observations)): # For every column
26        curr_observation = observations[i]
27        for key in scores: # For every row
28            if key in emission_matrix:
29                state_scores = []
30                for key2 in scores: # Finds all the cross ...
31                    probabilities
32                    previous_score = scores[key2][i-1]
33                    state_scores.append(
34                        previous_score*transition_matrix[key2][key]*emission_matrix[key][curr_observation])
35                scores[key].append(max(state_scores))
36            else:
37                scores[key].append(0.0)
38    # Adds the state to the prediction sequence
39    max_score = 0
40    best_key = ''
41    for key in scores:
42        score = scores[key][i]
43        if score > max_score:
44            max_score = score
45            best_key = key
46    prediction_seq.append(best_key)
47    for key in scores:
48        prob = scores[key][-1]
49        if prob > prediction_prob:
50            prediction_prob = prob
51    #print(scores)
52    return prediction_seq, prediction_prob
```



## 5 Transcript of a Sample Session

We wrote a test case to demonstrate our code. In the test case, the transition probabilities, emission probabilities and the initial probabilities are pre computed.

```
1 # TEST CASE WHETHER THE CODE WORKS
2 observe = ['a', 'c', '-']
3 def RunPredict(obs):
4     res1= forward(obs, emission_prob, transition_prob, initial_prob)
5     print("Forward Algorithm")
6     print("The transition matrix: ", res1[0])
7     print("The predicted Sequence: ", res1[1])
8     res2 = viterbi(transition_prob, emission_prob, initial_prob, obs)
9     print("Viterbi Algorithm")
10    print("The transition matrix: ", res2[0])
11    print("The predicted Sequence: ", res2[1])
12 RunPredict(observe)
```

The results are listed above.

```
1 Forward Algorithm
2 The transition matrix:
3 [[3.78930556e-04 3.46360108e-05 0.00000000e+00]
4  [5.16782099e-02 0.00000000e+00 3.84425840e-07]
5  [0.00000000e+00 9.85394754e-04 4.93823908e-07]
6  [0.00000000e+00 8.51664250e-06 2.86179325e-07]
7  [7.10353401e-04 0.00000000e+00 9.21598677e-07]
8  [5.08479464e-04 1.63574696e-05 6.58840382e-07]
9  [8.49887105e-04 0.00000000e+00 2.39314004e-06]
10 [0.00000000e+00 0.00000000e+00 3.75208499e-07]
11 [5.40837248e-04 0.00000000e+00 3.56272824e-07]
12 [0.00000000e+00 0.00000000e+00 1.38254914e-06]
13 [1.69977421e-04 1.10326874e-05 2.09844298e-06]
14 [5.06315722e-04 0.00000000e+00 3.93695485e-06]
15 [0.00000000e+00 1.25433842e-05 4.67463630e-08]
16 [1.01986453e-03 0.00000000e+00 2.31021295e-06]
17 [0.00000000e+00 5.44300096e-06 5.42705682e-07]
18 [0.00000000e+00 1.13609677e-05 2.68824171e-06]
19 [1.18984195e-03 6.71270853e-05 2.66194551e-06]
20 [6.99907027e-04 1.30036008e-05 2.45411322e-06]
21 [9.75280284e-04 0.00000000e+00 1.38623639e-07]
22 [0.00000000e+00 9.65695879e-07 0.00000000e+00]
23 [0.00000000e+00 1.32827758e-05 7.82009742e-08]]
24 The predicted Sequence: ['a', 'c', 'm']
25 Viterbi Algorithm
26 The transition matrix: ['a', 'c', 'm']
27 The predicted Sequence: 3.2981048769268155e-06
```

The predicted sequence using two algorithms are consistent.

## 6 Demos

### 6.1 Step 1

The client has to write observation for the program. The observation has to be an array, whose elements are either the standard amino acids or the deletion character "-". A more detailed abbreviations for amino acids are listed in Appendix section. Moreover, we only take lowercase letters as input.

### 6.2 Step 2

After the client decided a sequence of observations, for instance ["a", "c", "m"], "m"], pass the observation array into the function runPredict to see the result.

```
1 observation = ["a", "c", "m"]
2 RunPredict(observation)
```

## 7 Some Interesting Part

We have decided to put our programs in real life setting by using a test protein sequence. We used ferret's TRPA1 protein sequence and compared it with house mouse's TRPA1 protein sequence to get the aligned sequence (with deletion states). We then use forward algorithm to reconstruct the sequence. Finally we use EMBOSS Needle again to compare our reconstructed sequence with the original sequence to check whether it is accurate. We achieved 88.7 percent accuracy which is pretty good.

```
1 # Aligned_sequences: 2
2 # 1: EMBOSS_001
3 # 2: EMBOSS_001
4 # Matrix: EBLOSUM62
5 # Gap_penalty: 10.0
6 # Extend_penalty: 0.5
7 #
8 # Length: 1125
9 # Identity:      896/1125 (79.6%)
10 # Similarity:   998/1125 (88.7%)
11 # Gaps:         35/1125 ( 3.1%)
12 # Score: 4749.5
```

## 8 What We Have Learnt

### 8.1 Leon Zhang (leonz, 1664187)

1. How AI can be applied to Bioinformatics

2. Hidden Markov models
3. Viterbi Algorithms
4. Forward Algorithms
5. Team working skills

## 8.2 Qiubai Yu (qiubay, 1663777)

1. Bioinformatics
2. Probability theory
3. Dynamic programming
4. Team working skills

## 9 If We Had More Time

If we had more time, we wish to take more data samples by finding more mouse family animals. We also wish to take more time in comparing our results from using Viterbi and forward algorithms. Finally We think it would be cool to spend more time to compare our reconstructed sequence with the labeled sequence (the original sequence), so will be more confident that our project has practical significance.

## 10 Citations

1. Amino acid, [https://en.wikipedia.org/wiki/Amino\\_acid](https://en.wikipedia.org/wiki/Amino_acid)
2. Protein, <https://en.wikipedia.org/wiki/Protein>
3. 2009 Nature Structural and Molecular Biology, [www.translatome.net/significance/2009NSMB.html](http://www.translatome.net/significance/2009NSMB.html)
4. Pairwise sequence alignment tool, <https://www.ebi.ac.uk/Tools/psa/>
5. Forward Algorithm, <https://www.youtube.com/watch?v=M7afek1nEKM>

## 11 Partners' Reflections

Partners' names and roles are in section 2.

## 11.1 Leon Zhang's Reflection

Challenges:

1. Managing data files and cleaning them
2. Learning algorithms outside of class

Benefits:

1. learn about applications in unfamiliar area
2. Practice self learning skills

## 11.2 Qiubai Yu's Reflection

Challenges:

1. applications in unfamiliar area
2. programming skills like dynamic programming

Benefits:

1. learn about applications in unfamiliar area
2. learn about hidden markov model and probability theories

## 12 Appendix

Table of Standard Amino Acid Abbreviations <sup>[1]</sup>		
Amino acid	3-letter	1-letter
Alanine	Ala	A
Arginine	Arg	R
Asparagine	Asn	N
Aspartic acid	Asp	D
Cysteine	Cys	C
Glutamine	Gln	Q
Glutamic acid	Glu	E
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Leucine	Leu	L
Lysine	Lys	K
Methionine	Met	M
Phenylalanine	Phe	F
Proline	Pro	P
Serine	Ser	S
Threonine	Thr	T
Tryptophan	Trp	W
Tyrosine	Tyr	Y
Valine	Val	V