

CONSTRUCTION AND APPLICATION BASED ON COMPRESSING DEPICTION IN PROFILE HIDDEN MARKOV MODEL

Zhijian Zhou^{1,2}, Daoliang Li^{2,3}, Li Li^{2,3}, Zetian Fu^{2,3,*}

¹ College of Science, China Agriculture University, Beijing, China, 100083

² Key Laboratory of Modern Precision Agriculture System Integration, Research Ministry of Education, Beijing, China, 100083

³ College of Information and Electrical Engineering, China Agricultural University, Beijing, China, 100083

* Corresponding author, Address: P.O.Box 209#, China Agricultural University, East Campus, Beijing 100083, P. R. China, Tel:+86-10-62736323, Email:fzt@cau.edu.cn

Abstract: A method to express Profile Hidden Markov Model (Profile HMM) parameters with compressing matrix is presented, which is obtained by imposing the characteristics of both the state transfer and the character output in the Profile HMM.

Key words: Profile HMM, Multiple sequence alignment, Bioinformatics

1. INTRODUCTION

The Profile HMM is composed by a Markov chain including matching, insertion, delete states, and an observable stochastic process namely as observation chain. The state chain depicts the transfer relationship among different state. The observation chain represents the statistical association between the state and observation. Generally, the state in Markov process cannot be observed directly. It can only be understood through the observable process. The Profile HMM model was introduced to bioinformatics by Krogh (1994), and now was widely used in the Alignment of biological sequence.

A one grade Profile HMM with matching range as L is illustrated as Figure 1. This is a linear model. It progress only one direction from left to right. There are $3L+1$ system states in the system namely as matching (M), insertion (I), and delete (D) respectively. For convenience in coding the program, two extra states, say as beginning (M_0) and ending (M_{L+1}), are involved in the model. They do not output any character to influence the model activity. The character set depends on the concerning object. For example, there are four characters in DNA sequence denoted by A, G, C, and T, and twenty characters for amino acid sequence. In the Figure 1, rectangle, diamond, and circle denote matching, insertion and delete state respectively. The arrow connected different states indicates the state transfer relation and direction. In an ascertain model, the transfer probability and character release between different state is wholly determinate.

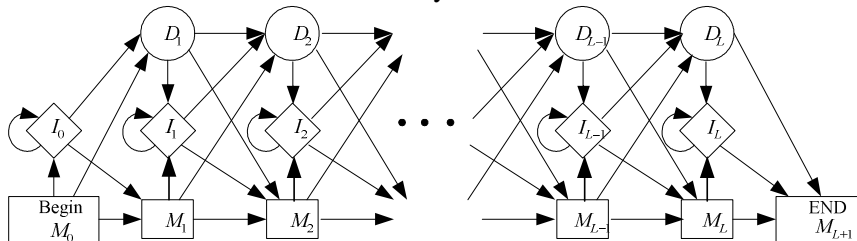


Figure 1 Model of one grade Profile HMM with matching range as L

2. COMPRESSION FORMS OF PROFILE HMM PARAMETERS

Based on the stochastic processes theory, a Profile HMM can be totally ascertained if $\lambda = \lambda(S, \Omega, A, B, \pi)$ is known. Where S denotes as state set, Ω is observation set, A is state transfer probability matrix, B is character output probability matrix, and π is initial state probability distribution function, respectively. Suppose a Profile HMM has matching range as L . Then the number of elements for S is $N = 3L + 1$. If two extra state, beginning and ending stats, are also included, the number of elements for S become as $N = 3(L + 1)$. So the order of state transfer matrix A is $3(L + 1) \times 3(L + 1)$. It can be found from Figure 1, at any state such as M_l , three front states as inputs of the state are at the most. They locate at same layer, say as $l-1$. This means only three occasions, $M_{l-1}M_l, I_{l-1}M_l$ and $D_{l-1}M_l$, can turn up when the state transfer from M_{l-1} to M_l . Similarly, only three occasions, M_lI_l, I_lI_l, D_lI_l and $M_{l-1}D_l, I_{l-1}D_l, D_{l-1}D_l$ for I_l and D_l state occur respectively. Remembering none of character output for delete state, the

state transfer probability matrix as order as $3(L+1) \times 3(L+1)$ can be compressed to $9 \times (L+1)$ without lost any information:

$$\bar{\mathbf{A}}_{9 \times (L+1)} = \begin{pmatrix} a[M_0M_1] & a[M_1M_2] & \cdots & a[M_{l-1}M_l] & \cdots & a[M_{L-1}M_L] & a[M_L M_{L+1}] \\ a[I_0M_1] & a[I_1M_2] & \cdots & a[I_{l-1}M_l] & \cdots & a[I_{L-1}M_L] & a[I_L M_{L+1}] \\ 0 & a[D_1M_2] & \cdots & a[D_{l-1}M_l] & \cdots & a[D_{L-1}M_L] & a[D_L M_{L+1}] \\ a[M_1I_1] & a[M_2I_2] & \cdots & a[M_lI_l] & \cdots & a[M_LI_L] & 0 \\ a[I_1I_1] & a[I_2I_2] & \cdots & a[I_lI_l] & \cdots & a[I_LI_L] & 0 \\ a[D_1I_1] & a[D_2I_2] & \cdots & a[D_lI_l] & \cdots & a[D_LI_L] & 0 \\ a[M_0D_1] & a[M_1D_2] & \cdots & a[M_{l-1}D_l] & \cdots & a[M_{L-1}D_L] & 0 \\ a[I_0D_1] & a[I_1D_2] & \cdots & a[I_{l-1}D_l] & \cdots & a[I_{L-1}D_L] & 0 \\ 0 & a[D_1D_2] & \cdots & a[D_{l-1}D_l] & \cdots & a[D_{L-1}D_L] & 0 \end{pmatrix}_{9 \times (L+1)}$$

Where $S = \{M_0, I_0, M_1, I_1, D_1, M_2, I_2, D_2, \dots, M_L, I_L, D_L, M_{L+1}\}$ are state set, and $a[X_i Y_j]$ represents the one step transfer probability from X_i to Y_j . The elements sign as “0” means no corresponding state transfer, and elements in l column denote the probability of one step state transfer ending at l layer. Let $\bar{\mathbf{A}}_{9 \times (L+1)} \equiv (\bar{a}_{ij})_{9 \times (L+1)}$, $\bar{a}_{ij} = a[X_i Y_j]$, $1 \leq i \leq 9$, $1 \leq j \leq L+1$, then following relations can be obtained ,

$$(1) i \bmod 3 = \begin{cases} 1, & X = M \\ 2, & X = I \\ 0, & X = D \end{cases} \quad (2-1)$$

$$(2) (d, Y) = \begin{cases} (j-1, M) & \text{if } 1 \leq i \leq 3 \\ (j, I) & \text{if } 4 \leq i \leq 6 \\ (j-1, D) & \text{if } 7 \leq i \leq 9 \end{cases} \quad (2-2)$$

$$(3) \bar{a}_{31} = \bar{a}_{91} = \bar{a}_{4(L+1)} = \cdots = \bar{a}_{9(L+1)} = 0 \quad (2-3)$$

$$(4) a_{ij} + a_{(i+3)(j-1)} + a_{(i+6)j} = 1, i = 1, 2, 3; \quad j = 2, 3, 4 \quad (2-4)$$

Similarly, let $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$ is observation character set, then character output probability matrix $\mathbf{B}_{K \times (3L+1)}$ can be compressed as

$$\bar{\mathbf{B}}_{K \times (2L+1)} = \begin{pmatrix} b_0^l[\omega_1] & b_1^M[\omega_1] & b_1^l[\omega_1] & b_2^M[\omega_1] & b_2^l[\omega_1] & \cdots & b_L^M[\omega_1] & b_L^l[\omega_1] \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ b_0^l[\omega_k] & b_1^M[\omega_k] & b_1^l[\omega_k] & b_2^M[\omega_k] & b_2^l[\omega_k] & \cdots & b_L^M[\omega_k] & b_L^l[\omega_k] \end{pmatrix}_{K \times (2L+1)} \quad (2-5)$$

Where $b_j^X[\omega_k]$ is the probability of output ω_k at the state X_j . Let $\bar{\mathbf{B}}_{K \times (2L+1)} \equiv (b_{ij})_{K \times (2L+1)}$, $b_{ij} = b_d^X(\omega_i)$, $1 \leq i \leq K$, $1 \leq j \leq 2L+1$, following relations can be obtained ,

$$j \bmod 2 = \begin{cases} 1, & (X, d) = (I, \frac{j-1}{2}) \\ 0, & (X, d) = (M, \frac{j}{2}) \end{cases} \quad (2-$$

6)

Therefore a profile HMM can be simplified by the compression state transfer probability matrix with order of $9 \times (L+1)$ and the compression character output probability matrix with order of $K \times (2L+1)$.

3. FORWARD ALGORITHM

Let us consider the observation sequence $O = (o_1, o_2, \dots, o_T)$. The matching range is L in Profile HMM (λ). Based on compress state transfer probability matrix $\bar{\mathbf{A}}_{9 \times (L+1)}$ and character output probability matrix $\bar{\mathbf{B}}_{K \times (2L+1)}$, forward algorithm can be obtained.

Definition 1 Let $\alpha_i^X(t) = P(o_1, o_2, \dots, o_t, \text{end of } X_t | \lambda)$, $X = M, I, D$ be the probability when part sequence $O_t = (o_1, o_2, \dots, o_t)$ output in X_t state at l ($1 \leq l \leq L$) layer. Then $(\alpha_i^M(t), \alpha_i^I(t), \alpha_i^D(t))^T$ is the probability vector for l layer, denoted as $\alpha_i(t)$.

Definition 2 Let $\varphi(X_r) = (a[M_{r-1}X_r], a[I_{r-1}X_r], a[D_{r-1}X_r])^T$, $X = M, D$ be a column vector composed by one step transfer probability from state X_{r-1} to X_r , and $\varphi(I_q) = (a[M_qI_q], a[I_qI_q], a[D_qI_q])^T$ also be a column vector composed by one step transfer probability from state I_{q-1} to I_q . Thus $\varphi(M_p) = (\bar{a}_{1p}, \bar{a}_{2p}, \bar{a}_{3p})^T$, $\varphi(I_q) = (\bar{a}_{4q}, \bar{a}_{5q}, \bar{a}_{6q})^T$ and $\varphi(D_r) = (\bar{a}_{7r}, \bar{a}_{8r}, \bar{a}_{9r})^T$, if compress state transfer probability matrix is $\bar{\mathbf{A}}_{9 \times (L+1)}$ and character output probability matrix is $\bar{\mathbf{B}}_{K \times (2L+1)}$.

Now the forward algorithm for Profile HMM can be express as,

1) Initiation

$$\begin{aligned} \alpha_0(0) &= (1, 0, 0)^T \\ \alpha_0(t) &= (0, 0, 0)^T, t = 1, 2, \dots, T+1 \end{aligned} \quad (3-1)$$

$$\alpha_l(0) = (0,0,0)^T, l = 1, 2, \dots, L+1$$

2) Recursion calculation

$$\alpha_l(t) = \begin{pmatrix} \alpha_{l-1}^T(t-1)\phi(M_l)b_l^M(o_t) \\ \alpha_l^T(t-1)\phi(I_l)b_l^I(o_t) \\ \alpha_{l-1}^T(t)\phi(D_l) \end{pmatrix} \quad t = 1, 2, \dots, T, \quad l = 1, 2, \dots, L \quad (3-$$

2)

where $b_d^M(o_t) = b_{i(2d)}, b_l^I(o_t) = b_{i(2d+1)}$, when $o_t = w_i$.

3) Ending

Thus

$$\alpha_{L+1}^M(T+1) = \alpha_L^M(T)a_{1(L+1)} + \alpha_L^I(T)a_{2(L+1)} + \alpha_L^D(T)a_{3(L+1)} \quad (3-$$

3)

and the probability is

$$P(O|\lambda) = \alpha_{L+1}^M(T+1) \quad (3-4)$$

4. APPLICATION OF PROFILE HMM

An example of compress state transfer probability matrix and character output probability matrix for known multiple sequence comparison is shown below. Suppose DAN sequences have be alimnt as Table 1. Matching states locate in first, second, and sixth column. Thus m is 5, T is 6 and L is 3. The corresponding Profile HMM framework is depicted in Figure 2.

Table 7 Alignment of multiple DNA sequence

	1	2	3	4	5	6
bat	A	G	-	-	-	C
rat	A	-	A	G	-	C
cat	A	G	-	A	A	-
gnat	-	-	A	A	A	C
goat	A	G	-	-	-	C

In order to avoid zero in probability calculation, pseudo counting is adapted. For example, character A appears 4 times in first column, the probability to output A for this state is $b_1^M(A) = \frac{4+1}{4+4} = 0.625$. Character G does not show up and the probability to output G is $b_1^M(G) = \frac{0+1}{4+4} = 0.125$.

Table 2 list character output probability at the state in which state transfer real occurs. As the same rule used, the compress state transfer probability matrix can be obtained as matrix (4-1).

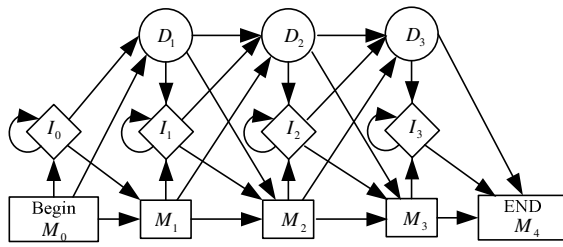


Figure 7 Profile HMM frameworks for the case of three states

Table 2. Frequency of character output

	$b_1^M(\omega_k)$	$b_2^M(\omega_k)$	$b_3^M(\omega_k)$	$b_4^M(\omega_k)$
A	0.625	0.143	0.333	0.125
G	0.125	0.571	0.333	0.125
C	0.125	0.143	0.111	0.625
T	0.125	0.143	0.222	0.125

$$A = \begin{pmatrix} 0.625 & 0.571 & 0.500 & 0.833 \\ 0.333 & 0.333 & 0.500 & 0.500 \\ 0.000 & 0.250 & 0.200 & 0.667 \\ 0.143 & 0.333 & 0.167 & 0.000 \\ 0.333 & 0.167 & 0.500 & 0.000 \\ 0.250 & 0.600 & 0.333 & 0.000 \\ 0.250 & 0.286 & 0.167 & 0.000 \\ 0.333 & 0.333 & 0.333 & 0.000 \\ 0.000 & 0.500 & 0.200 & 0.000 \end{pmatrix} \quad (4-1)$$

Suppose $O = AGC$ is a new observation sequence. Thus the probability that the sequence can be observed by using the model as shown above is $P(O|\lambda) \square 0.0373$

In practical application, the observation sequence can be classed according to its probability for different model λ_i that established by corresponding training data.

ACKNOWLEDGEMENTS

This work supported by the grant (10631070) from the National Natural Science Foundation of China.

REFERENCES

- Gong Guang-lu, Qian Min-Ping , 2004 , Application stochastic process course and stochastic model in intelligent calculation, Qinghua Publisher, Beijing.
- Kato A , Mian I S , Haussler D. , 1994 , A hidden Markov model that finds genes in E.coli DNA. *Nucleic Acids Research* , 22:4768-4778.
- Krogh A , Brown M , Mian I S , Sjolander K , and Haussler D. , 1994 , Hidden Markov models in computational biology: applications to protein modeling , *Journal of Molecular Biology* , 235:1501-1531.
- Wang Yu-Fei, Shi Ding-Hua , 2006 , Bioinformatic intelligent arithmetic and its application, Chemistry industry Publisher, Beijing.