

1 Pairwise HMM

Pairwise HMM とは 2 つの sequence data (DNA の塩基配列やタンパク質の residue 配列) の類似度をモデル化した HMM である . これは discrete HMM の拡張であると考えられるが , 一般的な discrete HMM とは出力が 2 つのシンボルの組になるという点で異なる . 隠れ状態の状態遷移図は以下になる . M は Match 状態 , X, Y はそれぞれの配列の挿入状態を表す .

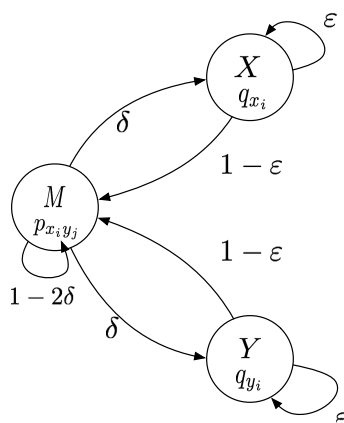


図 1 phmm の状態遷移図

1.1 Notations

定式化する上で Notation を整理する . 観測変数はシンボル列 x, y である .

$$\mathbf{x} = (x^1, \dots, y^{T_x}) \quad (1.1)$$

$$\mathbf{y} = (y^1, \dots, y^{T_y}) \quad (1.2)$$

$$(1.3)$$

観測変数 x^i, y^j に対応する隠れ変数 z^{ij} は以下のように定義する . ここで , x^i もしくは y^j の代わりに gap を許す (観測されたシンボルが対応しない場合がある) .

$$\mathbf{z} = \begin{pmatrix} z^{11} & \dots & z^{1T_y} \\ \vdots & \ddots & \vdots \\ z^{T_x 1} & \dots & z^{T_x T_y} \end{pmatrix} \quad (1.4)$$

隠れ変数は z^{ij} は状態数 K に対して 1-of- K 符号化方式で表される .

$$\mathbf{z}^{ij} = (z_1^{ij}, \dots, z_K^{ij}) \quad (1.5)$$

$$z_k^{ij} = \begin{cases} 1 & x^i, y^j \text{ は } k \text{ 番目の状態から出力される} \\ 0 & \text{otherwise} \end{cases} \quad (1.6)$$

また，表記の簡単のため， $z^{i0} = z^{0j} = 0$ ，遷移の方向の offset として $\delta = \{(\delta_x, \delta_y)\}$ と定義する．ここでは， $\delta = \{(1, 1), (0, 1), (1, 0)\}$ となる．これを用いて，隠れマルコフモデルの同時分布は以下のように表される．

$$p(\mathbf{x}, \mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) = p(\mathbf{z}^{11} | \boldsymbol{\alpha}) \prod_{i,j=1}^{T_x, T_y} \prod_{(\delta_x, \delta_y) \in \delta} p(\mathbf{z}^{ij} | \mathbf{z}^{i-\delta_x, j-\delta_y}, \boldsymbol{\beta}) \prod_{i,j=1}^{T_x, T_y} p(x^i y^j | \mathbf{z}^{ij}, \boldsymbol{\phi}) \quad (1.7)$$

上記で， $\boldsymbol{\theta} = (\boldsymbol{\alpha}, \boldsymbol{\beta}, \boldsymbol{\phi})$ は HMM のパラメタである．

$$\begin{aligned} \text{初期確率: } p(\mathbf{z}^{11} | \boldsymbol{\alpha}) &= \prod_{k=1}^K \alpha_k^{z_k^{11}} \\ \text{出力確率: } p(x^i, y^j | \mathbf{z}^{ij}, \boldsymbol{\phi}) &= \prod_{k=1}^K p(x^i, y^j | \phi_k)^{z_k^{ij}} \\ \text{遷移確率: } p(\mathbf{z}^{ij} | \mathbf{z}^{i-\delta_x, j-\delta_y}, \boldsymbol{\beta}) &= \prod_{l,k=1}^{L,K} \beta_{lk}^{z_l^{i-\delta_x, j-\delta_y} z_k^{ij}} \end{aligned} \quad (1.8)$$

ここで，

$$\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K) \in \mathbb{R}^K \quad (1.9)$$

$$\boldsymbol{\beta} = (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K) \in \mathbb{R}^K \times \mathbb{R}^K \text{ where } \boldsymbol{\beta}_k = (\beta_{k1}, \dots, \beta_{kK}) \in \mathbb{R}^K \quad (1.10)$$

$$\boldsymbol{\phi} = (\boldsymbol{\phi}_1, \dots, \boldsymbol{\phi}_K) \quad (1.11)$$

である．

1.2 EM algorithm

EM algorithm によりパラメタの最適化を行う．EM アルゴリズムは以下のようくり返しのアルゴリズムである．

$$\text{E-step} \quad (1.12)$$

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) = \mathbb{E}_{p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}^{old})} [\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})] \quad (1.13)$$

$$(1.14)$$

$$\text{M-step} \quad (1.15)$$

$$\boldsymbol{\theta}^{new} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) \quad (1.16)$$

Estep では，事後分布 $p(\mathbf{z}^{ij} | \mathbf{x})$ と $p(\mathbf{z}^{ij}, \mathbf{z}^{i-\delta_x, j-\delta_y} | \mathbf{x})$ を求める必要がある．この事後分布は forward-backward アルゴリズムで効率的に求めることが出来る．forward 変数 f と backward 変数 b を以下のように定義する．この2つの変数から，事後分布を求めることが出来る．ここでは，表記の簡単のためパラメタは明記しない．

$$f^{ij} = p(x^1, \dots, x^i, y^1, \dots, y^j, \mathbf{z}^{ij}) \quad (1.17)$$

$$b^{ij} = p(x^{i+1}, \dots, x^{T_x}, y^{j+1}, \dots, y^{T_y} | \mathbf{z}^{ij}) \quad (1.18)$$

$$\gamma^{ij} = p(\mathbf{z}^{ij} | \mathbf{x}, \mathbf{y}) = f^{ij} b^{ij} / p(\mathbf{x}, \mathbf{y}) \quad (1.19)$$

$$\gamma_k^{ij} = p(z_k^{ij} = 1 | \mathbf{x}, \mathbf{y}) \quad (1.20)$$

$$\xi^{ij\delta} = p(\mathbf{z}^{ij}, \mathbf{z}^{i-\delta_x, j-\delta_y} | \mathbf{x}, \mathbf{y}) = f^{i-\delta_x, j-\delta_y} p(\mathbf{z}^{ij} | \mathbf{z}^{i-\delta_x, j-\delta_y}) p(x^{ij} | \mathbf{z}^{ij}) b^{ij} / p(\mathbf{x}, \mathbf{y}) \quad (1.21)$$

$$\xi_{lk}^{ij\delta} = p(z_k^{ij} = 1, z_l^{i-\delta_x, j-\delta_y} = 1 | \mathbf{x}, \mathbf{y}) \quad (1.22)$$

ここで、表記の簡単のため γ, ξ を定義した。すると、これらは DP(Dynamic Programming) で計算できる。

$$f^{11} = p(z^{11})p(x^1, y^1 | z^{11}) \quad (1.23)$$

$$f^{ij} = p(x^i, y^j | z^{ij}) \sum_{(\delta_x, \delta_y) \in \delta} \sum_{z^{i-\delta_x, j-\delta_y}} p(z^{ij} | z^{i-\delta_x, j-\delta_y}) f^{i-\delta_x, j-\delta_y} \quad (1.24)$$

$$b^{T_x T_y} = \mathbf{1} \quad (1.25)$$

$$b^{ij} = \sum_{(\delta_x, \delta_y) \in \delta} \sum_{z^{i+\delta_x, j+\delta_y}} p(x^{i+\delta_x}, y^{j+\delta_y} | z^{i+\delta_x, j+\delta_y}) p(z^{i+\delta_x, j+\delta_y} | z^{ij}) b^{i+\delta_x, j+\delta_y} \quad (1.26)$$

Mstep では、この事後分布 $p(z|x, y, \theta^{old})$ を用いて Q 関数を最大化するパラメタ θ を求める。

$$Q(\theta, \theta^{old}) = \mathbb{E}_{p(z|x, y, \theta^{old})} [\ln p(x, y, z | \theta)] \quad (1.27)$$

$$\theta^{new} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^{old}) \quad (1.28)$$

それぞれのパラメタは以下の通り計算できる。

$$\alpha_k = \gamma_k^{11} / \sum_k \gamma_k^{11} \quad (1.29)$$

$$\beta_{lk} = \sum_{ij \in \delta} \xi_{lk}^{ij\delta} / \sum_{ij \in \delta} b \xi_{lk}^{ij\delta} \quad (1.30)$$

$$\theta_k^* = \underset{\theta_k}{\operatorname{argmax}} \sum_{ijk} \gamma_k^{ij} p(x_i, y_j | \theta_k, z_k^{ij} = 1) \quad (1.31)$$

θ は、emission の分布の形 $p(x, y | \theta, z)$ を仮定することで得ることが出来る。ここでは、シンボルを $s^1 \dots s^D$ の D 種類として、 s^m, s^n の組み合わせを出力する確率を μ^{mn} に持つ多項分布を仮定する。

$$p(x, y | \phi) = \prod_{mn} \mu_{mn}^{I(x, s^m) I(y, s^n)} \quad (1.32)$$

このとき、パラメータ $\phi = \mu$ はラグランジュの未定乗数法を用いて以下のように求められる。

$$\mu^{mn} = \sum_{ij} \gamma^{ij} I(x^i, s^m) I(y^j, s^n) / \sum_{ijmn} \gamma^{ij} I(x^i, s^m) I(y^j, s^n) \quad (1.33)$$

以上を用いて、パラメタの最適化を行うことができる。