

---

# HMM algorithm

---

## 1 Backward

$$P(Y_{k+1}, \dots, Y_n | X_k = C) = \sum_q P(P(Y_{k+2}, \dots, Y_n | t_{k+1} = q)P(q|C)P(x_{k+1}|q))$$

Define:

$$\beta_k(C) = P(Y_{k+1}, \dots, Y_n | X_k = C)$$

Inductive step:

$$\beta_k(C) = \sum_q \beta_{k+1}(q)P(q|C)P(Y_{k+1}|q)$$

## 2 Forward

$$P(Y_1, \dots, Y_n) = \sum_t (\prod_i P(Y_i | X_i)P(X_i | X_{i-1}))$$

$$P(Y_1, \dots, Y_k, X_k = q) = \sum_{q_1} P(Y_1, \dots, Y_k, X_{k-1} = q_1, X_k = q) = \sum_{q_1} P(Y_1, \dots, Y_{k-1}, X_{k-1} = q_1)P(X_k = q | X_{k-1} = q_1)P(Y_k | X_k = q)$$

Define:

$$\alpha_k(q) = P(Y_1, \dots, Y_k, X_k = q)$$

$$\alpha_1(q) = p(Y_1, X_1 = q) = P(X_1 = q | X_0)p(Y_1 | X_1 = q)$$

Inductive step:

$$\alpha_k(q) = \sum_{q_1} \alpha_{k-1}(q_1)P(X_k = q | X_{k-1} = q_1)P(Y_k | X_k = q)$$

## 3 Viterbi

$$\max_{X_1, \dots, X_k} P(X_1, \dots, X_k, Y_1, \dots, Y_k) = \max_q \max_{X_1, \dots, X_{k-1}} P(X_1, \dots, X_{k-1}, X_k = q, Y_1, \dots, Y_k)$$

Define:

$$\theta_k(q) = \max_{X_1, \dots, X_{k-1}} P(X_1, \dots, X_{k-1}, X_k = q, Y_1, \dots, Y_k)$$

Inductive step:

$$\theta_k(q) = \max_{q_1} \theta_{k-1}(q_1)P(X_k = q | X_{k-1} = q_1)P(Y_k | X_k = q)$$

## 4 Baum Welch

E step:

$$\gamma_t(j) = \alpha_t(j)\beta_t(j)/\alpha_T(q_F)$$

$$\xi_t(i, j) = \alpha_t(i)T_{ij}O_{j, O_{t+1}}\beta_{t+1}(j)/\alpha_T(q_F)$$

M step:

$$\hat{a}_{t,ij} = \sum_{t=1}^{T-1} \xi_t(i, j) / \sum_{t=1}^{T-1} \sum_{k=1}^N \xi_t(i, k)$$

$$\text{bhat}_j(v_k) = \sum_{t=1}^T \gamma_t(j) / \sum_{t=1}^T \xi_t(j)$$

## 5 Loss function for Baum Welch

EM optimize the log likelihood of the input:

$$\begin{aligned} & \log P(Y_1, \dots, Y_k | T, O, \pi) \\ &= \log \sum_{X_1, \dots, X_k} P(X_1, \dots, X_k, Y_1, \dots, Y_k | T, O, \pi) \\ &= \log \sum_{X_1, \dots, X_k} \prod_{i=1}^n P(X_i | X_{i-1}) P(Y_i | X_i) \end{aligned}$$

If I use  $\omega$  to represent  $Y_1, \dots, Y_k$ , and  $t$  to represent  $X_1, \dots, X_k$ , and  $\lambda$  to represent the model parameter.

Then according to Jensen's inequality,

$$\begin{aligned} & \log \sum_t P(\omega, t | \lambda) \\ &= \log \sum_t (P(\omega, t | \lambda) / P(t | \omega, \lambda^{(s)})) P(t | \omega, \lambda^{(s)}) \\ &\geq \sum_t P(t | \omega, \lambda^{(s)}) \log (P(\omega, t | \lambda) / P(t | \omega, \lambda^{(s)})) \end{aligned}$$

Define

$$\begin{aligned} f(\lambda) &= \log \sum_t P(\omega, t | \lambda) \\ g_s(\lambda) &= \sum_t P(t | \omega, \lambda^{(s)}) \log (P(\omega, t | \lambda) / P(t | \omega, \lambda^{(s)})) \end{aligned}$$

So we have  $f(\lambda) \geq g_s(\lambda)$

Optimizing  $g_s(\lambda)$

$$\begin{aligned} & g_s(\lambda) \\ &= \sum_t P(t | \omega, \lambda^{(s)}) \log P(\omega, t | \lambda) / P(t | \omega, \lambda^{(s)}) \\ &= \sum_t P(t | \omega, \lambda^{(s)}) (\log P(\omega, t | \lambda) - \log P(t | \omega, \lambda^{(s)})) \end{aligned}$$

$$\begin{aligned} & \max_{\lambda} g_s(\lambda) \\ &= \max_{\lambda} \sum_t P(t | \omega, \lambda^{(s)}) \log P(\omega, t | \lambda) \\ &= \max_{\lambda} \sum_t P(t | \omega, \lambda^{(s)}) \sum_i (\log P(t_i | t_{i-1}) + \log P(\omega_t | t_i)) \end{aligned}$$

In a word, we can minimize the negative log loss (or maximize the log likelihood  $\max_{\lambda} P(w | \lambda)$ ) by maximizing  $g_s(\lambda) = \sum_t P(t | w, \lambda^{(s)}) \log P(w, t | \lambda)$  at iteration  $s$ .

$\max_{\lambda} g_s(\lambda)$  has a closed-form solution such that the parameters can be directly calculated by the expected counts.