

$$p(x_t | x_{t-1}, x_{t-2}, \dots, x_1) = p(x_t | x_{t-1})$$

$$p(x_1 = i_1, \dots, x_T = i_T) = \pi(i_1) \cdot a_{i_1 i_2} \dots a_{i_{T-1} i_T}$$

$$D = \{ \bar{x}, \bar{y} \} \quad p(D | \lambda) \rightarrow \max_{\lambda = (\pi, A)}$$



$$x_t \in \{1, \dots, n\} \\ p(x_t = i | x_{t+1} = j) = a_{ji}$$

$$A = (a_{ij})_{i,j=1}^n$$

$p(\pi)$
 $p(A)$

$$\pi_i^{ML} = p(x_1 = i) = \frac{\# \{ q_{n1} = i \}}{N}$$

$$a_{ij}^{ML} = \frac{\# \{ q_{nt} = i, q_{n,t+1} = j \}}{\# \{ q_{nt} = i, q_{n,t+1} = x \}}$$

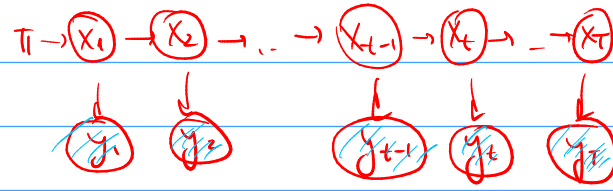
$$b_i^{ML}(k) = \frac{\# \{ q_{nt} = i, d_{nt} = k \}}{\# \{ q_{nt} = i \}}$$

$$\pi: p(x_1 = i) = \pi_i \quad p(Q, D | \lambda) = p(\bar{x} = q_1, q_2, \dots, q_T, \bar{y} = d_1, d_2, \dots, d_T | \pi, A, B) = \pi_{q_1} b_{q_1}(d_1) a_{q_1 q_2} b_{q_2}(d_2) a_{q_2 q_3} \dots a_{q_{T-1} q_T} b_{q_T}(d_T)$$

$$A: p(x_t = j | x_{t-1} = i) = a_{ij}$$

$$B: p(y_t = k | x_t = i) = b_i(k)$$

$$D = \{ d_{n1}, d_{n2}, \dots, d_{nT} \}_{n=1}^N$$

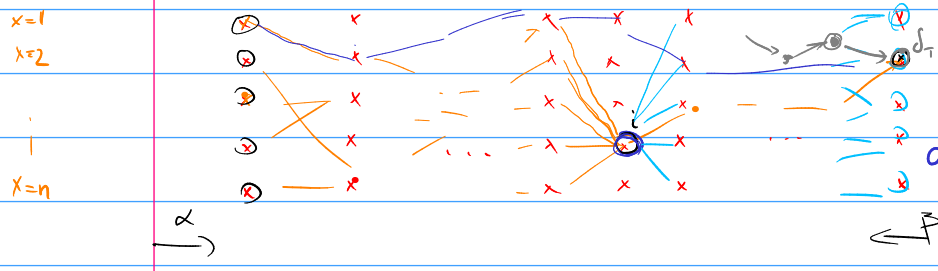


$$p(D | \lambda) = \sum_Q p(Q, D | \lambda) = \sum_{q_1, \dots, q_T} \pi_{q_1} \dots b_{q_T}(d_T)$$



$$p(D | \lambda) = \sum_{q_1, \dots, q_T} \pi_{q_1} \dots b_{q_T}(d_T)$$

- ① $p(D | \lambda) = ?$ *inference*
- ② $Q^* = \arg \max_Q p(Q, D | \lambda)$
- ③ $\lambda^* = \arg \max_{\lambda} p(D | \lambda)$ *learning*



$$d_t(i) = p(d_1, d_2, \dots, d_t, q_t = i | \lambda)$$

$$d_t(i) = \sum_{j=1}^n \alpha_{t-1}(j) a_{ji} b_i(d_t)$$

$$d_T(i) = p(D, q_T = i | \lambda)$$

$$p(D | \lambda) = \sum d_T(i)$$

$$\beta_t(i) = p(d_{t+1}, \dots, d_T | q_t = i, \lambda)$$

$$\beta_t(i) = p(d_{t+1}, \dots, d_T | q_t = i, \lambda) = \sum_{j=1}^n p(d_{t+1}, \dots, d_T | q_{t+1} = j, q_t = i, \lambda) = \sum_{j=1}^n p(d_{t+2}, \dots, d_T | q_{t+1} = j, q_t = i, \lambda) \cdot p(d_{t+1} | q_{t+1} = j, q_t = i, \lambda) \cdot p(q_{t+1} = j | q_t = i, \lambda)$$

$$\alpha_t(i) = p(d_1, \dots, d_t, q_t = i | \lambda) = \pi_i \cdot b_i(d_1)$$

$$\alpha_t(i) = p(d_{1:t}, q_t = i | \lambda) = p(d_{1:t-1}, d_t, q_t = i | \lambda) = \sum_{j=1}^n p(d_{1:t-1}, d_t, q_t = i, q_{t-1} = j | \lambda) = \sum_{j=1}^n \underbrace{p(d_{1:t-1}, q_{t-1} = j | \lambda)}_{\alpha_{t-1}(j)} \cdot \underbrace{p(q_t = i | q_{t-1} = j, \lambda)}_{a_{ji}} \cdot \underbrace{p(d_t | q_t = i, \lambda)}_{b_i(d_t)}$$

$$\beta_t(i) = \sum_{j=1}^{t+1} a_{ij} \cdot b_j(d_{t+1}) \cdot \beta_{t+1}(j)$$

$$\beta_1(i) = p(d_2, \dots, d_T | q_1 = i, \lambda)$$

$$p(D | \lambda) = \sum_{i=1}^n \pi_i b_i(d_1) \cdot \beta_1(i)$$

② $Q^* = \underset{Q}{\operatorname{argmax}} p(Q|D, \lambda)$
 $\rightarrow \underline{q_t^*} = \underset{q_t}{\operatorname{argmax}} p(q_t|D, \lambda)$

$\gamma_t(i) = p(q_t=i|D, \lambda) = \frac{p(q_t=i, d_{1:t}, d_{t+1:T}|\lambda)}{p(d_{1:t}, q_t=i|\lambda) \cdot p(d_{t+1:T}|q_t=i, \lambda)} = \frac{\alpha_t(i) \beta_t(i)}{p(D|\lambda)}$

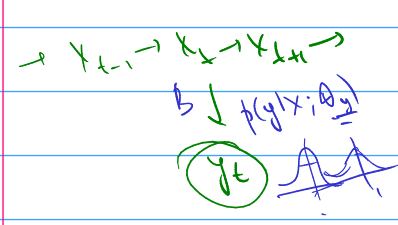
$\delta_t(i) = \max_{q_1, \dots, q_{t-1}} p(q_1, \dots, q_{t-1}, q_t=i, d_1, \dots, d_t|\lambda) = \max_{q_1, \dots, q_{t-1}} p(q_1, \dots, q_{t-2}, q_{t-1}, d_1, \dots, d_{t-1}|\lambda) \cdot \underbrace{p(q_t=i|q_{t-1}, d_1, \dots, d_{t-1}, \lambda)}_{a_{q_{t-1}, i}} \cdot \underbrace{p(d_t|q_t, \lambda)}_{b_i(d_t)}$

$\delta_t(i) = \max_j [\delta_{t-1}(j) \cdot a_{j, i} \cdot b_i(d_t)]$

③ Baum-Welch algorithm

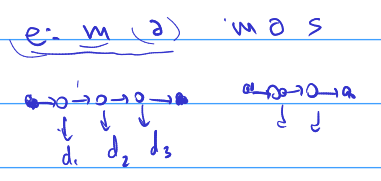
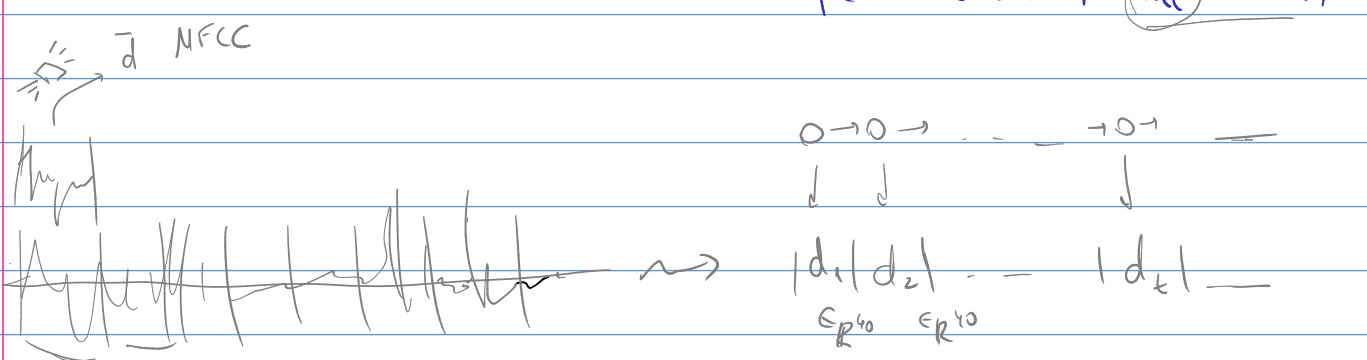
$E\text{-step } E[Q|\lambda^{(m)}] \leftarrow \left[\delta_t(i), \beta_t(i), \gamma_t(i), \varepsilon_t(i, j) = p(q_{t-1}=i, q_t=j|D, \lambda) \right]$
 $M\text{-step } \lambda^{(m+1)} := \operatorname{argmax}_{\lambda^{(m)}} E \log p(D, Q|\lambda)$

$\pi_i := \frac{E \# \{q_t=i\}}{N} = \frac{\sum_n \delta_t^{(m)}(i)}{N}$
 $a_{ij} = \frac{E \# \{q_t=i, q_{t+1}=j\}}{E \# \{q_t=i\}} = \frac{\sum_n \sum_{t=1}^{T-1} \varepsilon_{t+1}(i, j)}{\sum_n \sum_{t=1}^{T-1} \delta_t^{(m)}(i)}$
 $b_i(k) = \frac{E \# \{q_t=i, d_t=k\}}{E \# \{q_t=i\}} = \frac{\sum_n \sum_{t: d_t=k} \delta_t^{(m)}(i)}{\sum_n \sum_{t=1}^T \delta_t^{(m)}(i)}$
 $\delta_t^{(m)}(i) = p(q_t=i|d_{1:t}, \lambda)$



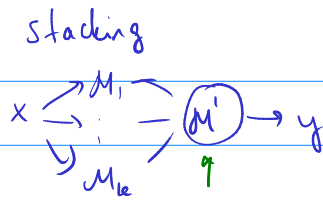
$\theta_y = \operatorname{argmax} \prod p(d_t|\dots)$
 $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$

$x_t=1 \rightarrow x_{t+1}=1 \rightarrow x_{t+2}=1$
 $p(z \text{ words } | x=i) = \prod_{i=1}^T (1 - a_{ii})$



$$p(M|D) \propto p(D|M) p(M)$$

M_1, M_2, \dots model combination



blending

$$y \sim \sum \alpha_k M_k(x)$$

Boosting

$$y \sim F_m(x) = \alpha_1 f_1(x) + \alpha_2 f_2(x) + \dots + \alpha_m f_m(x)$$

$$f_m(x), \alpha_m \sim F_{m-1}(x)$$

$$F_m(x) = \sum_{k=1}^m \alpha_k f_k(x)$$

AdaBoost

XGBoost - gradient boosting

