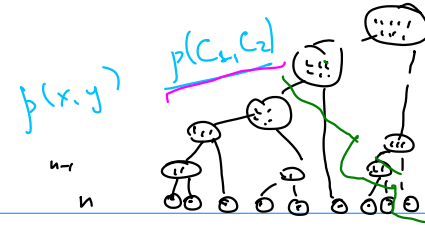


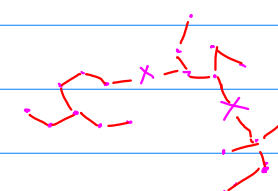
Clustering

$p(\bar{x}, \bar{y})$
 $[ACGTGG]$
 $[AA-GTGG]$

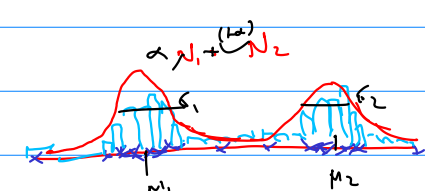
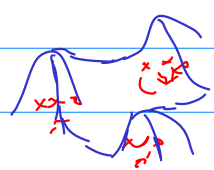
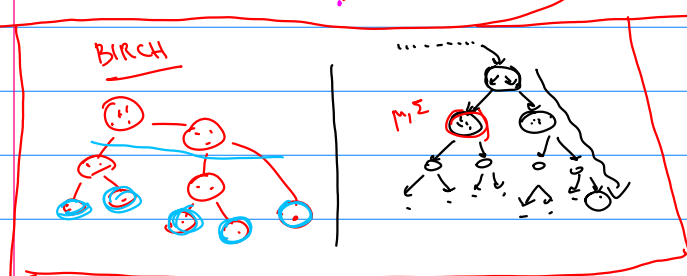
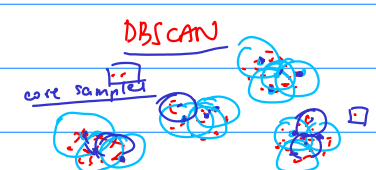


$p(c_1, c_2) =$
 single-link $\rightarrow \min_{x \in C_1, y \in C_2} p(x, y)$
 complete-link $\rightarrow \max_{x \in C_1, y \in C_2} p(x, y)$
 avg $p(x, y)$

\bar{x}_i, \bar{x}_n
 $W_{ij} = p(\bar{x}_i, \bar{x}_j)$



MST
 minimal spanning tree

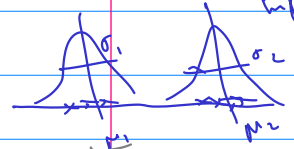


$$p(x | \mu_1, \mu_2, \sigma_1, \sigma_2, d) = d \cdot \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{1}{2\sigma_1^2}(x-\mu_1)^2} + (1-d) \cdot \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2\sigma_2^2}(x-\mu_2)^2}$$

 $\theta = \{d, \mu_1, \mu_2, \sigma_1, \sigma_2\}$
 $\theta_{ML} = \arg \max_{\theta} p(x | \theta)$
 $p(\theta | x) = \prod_{x \in D} [d \cdot e^{-\dots} + (1-d) \cdot e^{-\dots}] \rightarrow \max?$

Process: $x_n, z_n = \begin{cases} 1, & x_n \in C_1 \\ 0, & x_n \in C_2 \end{cases}$ $D = (x_n, z_n)$

$p(x, z | \theta) = \prod_{(x_n, z_n)} [d \mathcal{N}(x_n | \mu_1, \sigma_1)]^{z_n} [(1-d) \mathcal{N}(x_n | \mu_2, \sigma_2)]^{1-z_n}$



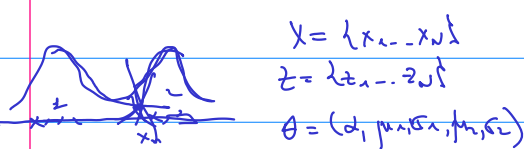
$$L_p(x, z | \theta) = \sum_n [z_n \ln(d \mathcal{N}_1) + (1-z_n) \ln((1-d) \mathcal{N}_2)] \rightarrow \max$$

$$= \sum_n [z_n \ln d + z_n \ln \mathcal{N}_1 + (1-z_n) \ln(1-d) + (1-z_n) \ln \mathcal{N}_2] \rightarrow \max$$

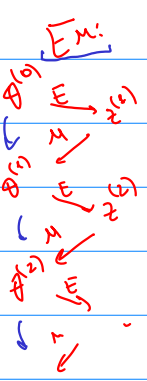
$\mu_1^*, \sigma_1^* : \sum_n z_n \ln \mathcal{N}(x_n | \mu_1, \sigma_1) = \sum_n \ln \mathcal{N}(x_n | \mu_1, \sigma_1) \rightarrow \max$
 $\mu_2^*, \sigma_2^* : \sum_n (1-z_n) \ln \mathcal{N}(x_n | \mu_2, \sigma_2) = \sum_n \ln \mathcal{N}(x_n | \mu_2, \sigma_2) \rightarrow \max$
 $\mu_1 = \frac{1}{d} \sum_{z_n=1} x_n, \sigma_1 = \frac{1}{d} \sqrt{\sum_{z_n=1} (x_n - \mu_1)^2}$

$\alpha^* = \frac{|C_1|}{N}$
 $\sum_n [z_n \ln d + (1-z_n) \ln(1-d)]$
 $|C_1| \ln d + |C_2| \ln(1-d)$
 $\frac{|C_1|}{d} - \frac{|C_1|}{1-d} = 0$
 $\alpha = \frac{|C_1|}{|C_1| + |C_2|}$

X $p(x | \theta) \rightarrow \max$ EM - expectation-maximization
 X, z $p(x, z | \theta) \rightarrow \max$ EM - expectation-maximization
 E-step - $z^{(n)} = \mathbb{E}_{p(z | \theta^{(n)}, x)}$
 M-step - $\theta^{(n+1)} = \arg \max_{\theta} p(x, z^{(n)})$

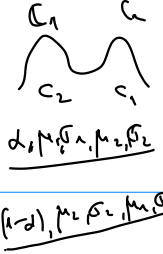


$X = \{x_1, \dots, x_n\}$
 $Z = \{z_1, \dots, z_n\}$
 $\theta = (d, \mu_1, \sigma_1, \mu_2, \sigma_2)$
 $p(x | \theta) = \prod (d \cdot \mathcal{N}_1 + (1-d) \cdot \mathcal{N}_2) \rightarrow \max$
 $p(x, z | \theta) = \prod (d \cdot \mathcal{N}_1)^{z_n} ((1-d) \cdot \mathcal{N}_2)^{1-z_n}$



$\alpha^{(n+1)} = \frac{\sum z_n}{N}$
 $\mu_1^{(n+1)} = \frac{1}{\sum z_n} \cdot \sum z_n x_n$
 $\sigma_1^{(n+1)} = \sqrt{\frac{1}{\sum z_n} \sum z_n (x_n - \mu_1^{(n+1)})^2}$

- init $\theta^{(0)}$ randomly
 - E-step: $\mathbb{E} z_n = p(x_n \in C_1 | \theta^{(n)}) = \frac{p(c_1) p(x_n | c_1)}{p(c_1) p(x_n | c_1) + p(c_2) p(x_n | c_2)}$
 $\mathbb{E} z_n = \frac{d \mathcal{N}(x_n | \mu_1, \sigma_1)}{d \mathcal{N}(x_n | \mu_1, \sigma_1) + (1-d) \mathcal{N}(x_n | \mu_2, \sigma_2)}$
 - M-step: $p(x, z^{(n)} | \theta) = \prod_n (d \cdot \mathcal{N}(x_n | \mu_1, \sigma_1))^{z_n} ((1-d) \mathcal{N}(x_n | \mu_2, \sigma_2))^{1-z_n}$
 $L_{mp} = \sum_n [z_n \ln d + (1-z_n) \ln(1-d) + z_n \ln \mathcal{N}(x_n | \mu_1, \sigma_1) + (1-z_n) \ln \mathcal{N}(x_n | \mu_2, \sigma_2)] \rightarrow \max$

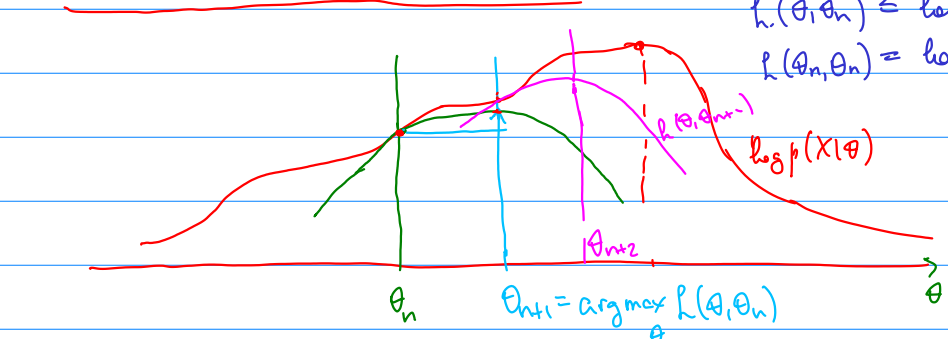


$(x_n, z_n) \quad p(x, z | \theta) \rightarrow \max$
 $p(x | \theta) = \int p(x, z | \theta) dz \rightarrow \max$
 $\log p(x | \theta) = \sum_n [z_n \log N_1 + (1 - z_n) \log N_2] \rightarrow \max$
 $Q = E_{z|x, \theta_n} \log p(x, z | \theta)$
 $Q(\theta_n) = \sum_k [E_{z|x, \theta_n} z_k \log N_1 + (1 - E_{z|x, \theta_n} z_k) \log N_2]$
 $\frac{d_n N(x_k | \mu_1, \sigma_1^2)}{d_n N(x_k | \mu_1, \sigma_1^2) + (1 - d_n) N(x_k | \mu_2, \sigma_2^2)}$

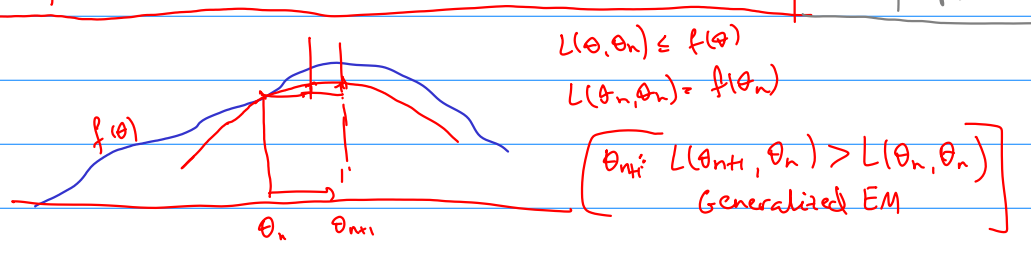
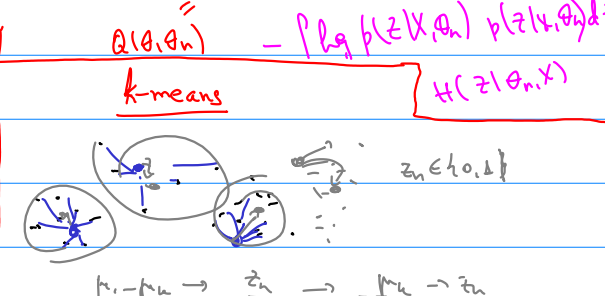
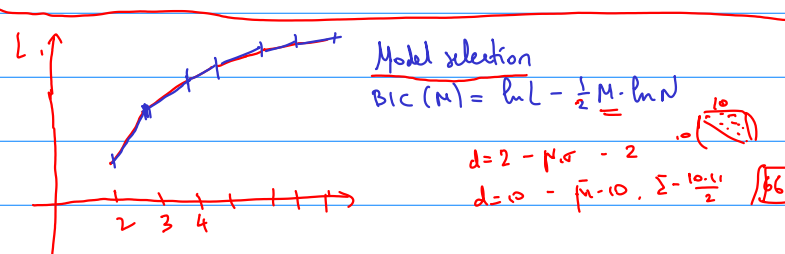
EM:
 $Q(\theta, \theta_n) = E_{p(z|x, \theta_n)} \log p(x, z | \theta)$
 $\theta_{n+1} := \arg \max_{\theta} Q(\theta, \theta_n)$
 $\text{Yes: } \log p(x | \theta_n) < \log p(x | \theta_{n+1})$
 $\log p(x | \theta) - \log p(x | \theta_n) =$
 $= \log \int p(x, z | \theta) dz - \log p(x | \theta_n) =$
 $= \log \int p(x, z | \theta) \cdot \frac{p(z | x, \theta_n)}{p(z | x, \theta_n)} dz - \log p(x | \theta_n)$
 $= \log \int \frac{p(x, z | \theta)}{p(z | x, \theta_n)} \cdot p(z | x, \theta_n) dz - \log p(x | \theta_n) \geq$

$\int f(z) p(z | x, \theta_n) dz = E_{z|x, \theta_n} f(z)$
 $\log E_{z|x, \theta_n} f(z) \geq E_{z|x, \theta_n} \log f(z)$ Jensen's inequality
 $g(\alpha x + (1 - \alpha)y) \geq \alpha g(x) + (1 - \alpha)g(y)$
 $g(\sum d_k x_k) \geq \sum d_k g(x_k), \sum d_k = 1$
 $g(E x) \geq E g(x)$

$\geq \int \log \frac{p(x, z | \theta)}{p(z | x, \theta_n)} p(z | x, \theta_n) dz - \log p(x | \theta_n) =$
 $= \int \log \frac{p(z | x, \theta) p(x | \theta)}{p(z | x, \theta_n) p(x | \theta_n)} p(z | x, \theta_n) dz$



$L(\theta, \theta_n) \leq \log p(x | \theta)$
 $L(\theta_n, \theta_n) = \log p(x | \theta_n)$
 $\int \log \frac{p(x, z | \theta)}{p(z | x, \theta_n)} p(z | x, \theta_n) dz$
 $= \int \log p(x, z | \theta) p(z | x, \theta_n) dz -$



$E_z f(z) \approx \frac{1}{R} \sum_{r=1}^R f(z_r)$
 $z_r \sim p(z | x, \theta_n)$

$\theta_{n+1} = \arg \max_{\theta} E_{z|x, \theta_n} \log p(x, z | \theta)$
Stochastic EM
 $= \int \log p(x, z | \theta) p(z | x, \theta_n) dz \approx \frac{1}{R} \sum_{r=1}^R \log p(x, z_r | \theta)$
 $\theta \rightarrow \max$

C_1
 $p_{1,0}$
 x_i
 x_m

m
 $ACGTACGT$
 $AGGTCGGT$
 $AGCTCCGT$
 $CGAACGGA$

$f(x, y) = \frac{\# \{x_i = y_i\}}{|x| \times |y|}$
 $C_j: \triangle \triangle \dots \triangle m$ (3m)
 P_{1A} P_{1C} ... P_{1m}
 P_{2A} P_{2C} ... P_{2m}
 P_{3A} P_{3C} ... P_{3m}

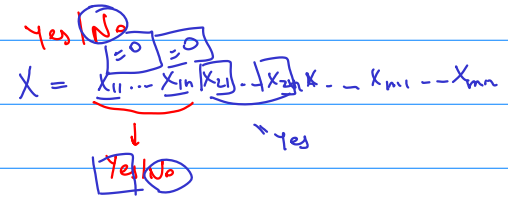
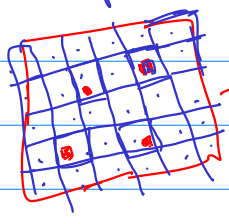
- init θ
 $C_1 \dots C_m$
 $E z_{nk} = \frac{\alpha_k p(x_n | C_k)}{\sum_j p(x_n | C_j)}$
 $\alpha_k = \frac{\# \{ \text{sykb } l \text{ ha nej } m \in C_k \}}{\# \{ \text{sykb } l \text{ ha nej } m \in C_1 \} + \dots + \# \{ \text{sykb } l \text{ ha nej } m \in C_m \}} = \frac{|C_k|}{|C_1| + \dots + |C_m|}$

$E_{z|x, p} p(x, z | p_{kml}) = E_z \log \prod_n \prod_k [\dots]^{z_{nk}}$

$p_{kml} := \frac{\sum E z_{nk}}{n \times x_{m=l}}$
 $\sum E z_{nk}$

$z \in \left[\sum_n z_{nk} \cdot \sum_m [\dots] \right] \rightarrow \max$
 $\sum E z_{nk} \cdot [\dots] \rightarrow \max$

Multiple instance learning (MIL)



x_{ij} z_{ij}
 $z_{ij} := \frac{p(x_{ij} = \text{Yes})}{\Sigma}$