

Ex. Mixture of (univariate) normal distributions

$$x_1, \dots, x_m \sim (1-\pi) N(\mu_0, \sigma_0^2) \oplus \pi N(\mu_1, \sigma_1^2)$$

$$\Theta = (\pi, \mu_0, \mu_1, \sigma_0^2, \sigma_1^2) \text{ mixture}$$

observed likelihood:

$$L_{\text{obs}, m}(\Theta) = \prod_{i=1}^m \left[(1-\pi) \varphi_0(x_i) + \pi \varphi_1(x_i) \right], \quad \varphi_j(x) = \varPhi\left(\frac{x-\mu_j}{\sigma_j}\right) \frac{1}{\sigma_j}$$

We know that for (e.g.) $\mu_0 = x_1$ and $\sigma_0 \rightarrow 0$, $L_{\text{obs}, m} \rightarrow \infty$ but there still exists a solution to the MLE equations that is a consistent estimator of the true Θ .

Numerical optimization is not a good method:

- L has singularities ($L \rightarrow \infty$)
- the solution we search for is not unique (symmetry of the problem; i.e. $L(\hat{\pi}, \hat{\mu}_0, \hat{\mu}_1, \hat{\sigma}_0^2, \hat{\sigma}_1^2) = L(1-\hat{\pi}, \hat{\mu}_1, \hat{\mu}_0, \hat{\sigma}_1^2, \hat{\sigma}_0^2)$)
and there will be saddle points where $\nabla L = 0$ but we might not have a root of $L' = 0$ an extreme
- $\sigma_j > 0$ (or Σ_j pos. def. may be difficult to enforce)
- the number of parameters is large (for d-variate distributions and q mixing components $d\% + q(d(d-1)/2) = O(d^2\%)$)

EM-algorithm: augmented dataset This is how we

$$(x_1, z_1), \dots, (x_m, z_m) \text{ iid}, \quad z_i \sim \text{alt}(\pi) \quad \leftarrow \text{generate from } x_i$$

$$x_i | z_i \sim \begin{cases} N(\mu_0, \sigma_0^2) & \text{if } z_i = 0 \\ N(\mu_1, \sigma_1^2) & \text{if } z_i = 1 \end{cases}$$

 (x_i, z_i) has a density w.r.t. $\lambda \otimes$ counting measure δ_{z_i} .

$$f_{x, z}(x, z) = f_{x|z}(x|z) P(z=z) = \begin{cases} \varphi_0(x)(1-\pi) & \text{if } z=0 \\ \varphi_1(x)\pi & \text{if } z=1 \end{cases} =$$

$$- [\varphi_0(x)(1-\pi)]^{1-\pi} [\varphi_1(x)\pi]^\pi \quad x \in \mathbb{R}, \pi \in \{0, 1\}$$

$$f_x(x) = \sum_{n=0}^1 f_{x, z}(x, z) = \varphi_0(x)(1-\pi) + \varphi_1(x)\pi = \underline{L_{\text{obs}, 1}(\Theta)}$$

z_i are unobserved. Complete likelihood

$$L_m^c(\theta) = \prod_{i=1}^m \left[\varphi_o(x_i)(1-\pi) \right]^{1-\tilde{z}_i} \left[\varphi_1(x_i)\pi \right]^{\tilde{z}_i}$$

$$\ell_m^c(\theta) = \sum (1-\tilde{z}_i) \log \varphi_o(x_i) + \sum (1-\tilde{z}_i) \cdot \log(1-\pi) \\ + \sum z_i \log \varphi_1(x_i) + \sum z_i \log \pi$$

E-step: for $\tilde{\theta}$ an initial estimate given

$$E_{\tilde{\theta}}[\ell_m^c(\theta) | x_1, \dots, x_m] = " \ell_m^c(\theta) \text{ with } \tilde{z}_i \text{ in place of } z_i "$$

$$E_{\tilde{\theta}}[\tilde{z}_i | x_1, \dots, x_m] = E_{\tilde{\theta}}[z_i | x_i] = P_{\tilde{\theta}}(z_i=1 | x_i) =$$

$$= \frac{f_{x|z_i}(x_i | 1) P(z_i=1)}{P f_x(x_i)} = \frac{\tilde{\varphi}_1(x_i) \tilde{\pi}}{\tilde{\varphi}_0(x_i)(1-\tilde{\pi}) + \tilde{\varphi}_1(x_i)\tilde{\pi}} =: \tilde{z}_i$$

for $\varphi_{\tilde{\theta}}(x) = \frac{1}{\tilde{\sigma}_{\tilde{\theta}}} \varphi\left(\frac{x - \tilde{\mu}_{\tilde{\theta}}}{\tilde{\sigma}_{\tilde{\theta}}}\right)$ with the initial estimated $\tilde{\theta}$ in place of θ .

M-step: $Q(\theta, \tilde{\theta}) := E_{\tilde{\theta}}[\ell_m^c(\theta) | x_1, \dots, x_m]$

$$\frac{\partial Q}{\partial \pi} = -\frac{\sum(1-\tilde{z}_i)}{1-\pi} + \frac{\sum \tilde{z}_i}{\pi} \stackrel{!}{=} 0 \Rightarrow \hat{\pi} = \frac{1}{m} \sum_{i=1}^m \tilde{z}_i$$

Binom. (Bernoulli)

$$\frac{\partial Q}{\partial \mu_0} = -\sum(1-\tilde{z}_i) \frac{(x_i - \mu_0)}{2\sigma_0^2} \stackrel{!}{=} 0 \Rightarrow \hat{\mu}_0 = \frac{\sum(1-\tilde{z}_i)x_i}{\sum(1-\tilde{z}_i)}$$

$$\frac{\partial Q}{\partial \sigma_0^2} = \sum(1-\tilde{z}_i) \left[-\frac{1}{2\sigma_0^2} + \frac{(x_i - \mu_0)^2}{2\sigma_0^4} \right] \stackrel{!}{=} 0 \Rightarrow \hat{\sigma}_0^2 = \frac{\sum(1-\tilde{z}_i)(x_i - \hat{\mu}_0)^2}{\sum(1-\tilde{z}_i)}$$

$\hat{\mu}_1$ and $\hat{\sigma}_1^2$ analogously.

Now set new $\tilde{\theta}$ to be $\hat{\theta}$ from the M-step, and repeat until convergence.

0. $\tilde{\theta} = \text{mt.}$
1. $\tilde{z}_i := \tilde{\varphi}_1(x_i)\tilde{\pi} / (\tilde{\varphi}_0(x_i)(1-\tilde{\pi}) + \tilde{\varphi}_1(x_i)\tilde{\pi})$
2. $\hat{\pi} = \frac{1}{m} \sum \tilde{z}_i$, $\hat{\mu}_0 = \frac{\sum(1-\tilde{z}_i)x_i}{\sum(1-\tilde{z}_i)}$, $\hat{\sigma}_0^2 = \dots$
3. $\tilde{\theta} := \hat{\theta}$ and go to step 1.

Week 11 - Missing data, EM algorithm

$(x_i, y_i) \sim N(\mu, \Sigma)$ random sample
 of size n

x_i missing for $i = m_0 + 1, \dots, m_0 + m_1$
 y_i missing for $i = m_0 + m_1 + 1, \dots, n$

CCA: complete case analysis - consider only i such that both x_i and y_i are known

ACA: available case analysis - for each estimator use all available observations

Simple Imputation: if x_i is known but y_i not, build model $y \sim \beta_0 + \beta_1 x_i$
 based on fully observed data and estimate y_i by $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \sim E[y_i | x_i]$

Parametric Imputation: estimate also ε_i in $y = \beta_0 + \beta_1 x_i + \varepsilon_i$ and
 generate resamples of $\hat{\varepsilon}_i \rightarrow \hat{y}_i := \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\varepsilon}_i$

Multiple Imputation: Repeat an imputation procedure M times, estimate $\hat{\theta}_j$
 in the j -th imputation, and not $\hat{\theta}_{MI} = \hat{\theta}_M \sim$ Bootstrap - get parameter
 estimates out of resamples \rightarrow possible to approximate also the variance
 of $\hat{\theta}_{MI}$: $\text{var}(\hat{\theta}_{MI}) = E(\text{var}(\hat{\theta}_{MI} | \text{imputed data})) + \text{var}(E(\hat{\theta}_{MI} | \text{imp. d}))$
 $\approx \frac{1}{M} \sum_{j=1}^M \text{var}(\hat{\theta}_{MI} | \text{imp. d}_j) + \frac{1}{M} \sum_{j=1}^M (\hat{\theta}_j - \hat{\theta}_{MI})(\hat{\theta}_j - \hat{\theta}_{MI})^\top$

but this is not a good estimator - needs to be adjusted for estimation of
 the sampling distribution

$$m_1 = m - m_0 - m_2 \quad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \quad r_i = \begin{pmatrix} x_i \\ y_i \end{pmatrix}$$

EM-algorithm: complete log-likelihood: $\theta = (\mu_1, \mu_2, \sigma_{11}, \sigma_{12}, \sigma_{22})^\top$

$$l_m^c(\theta) = c - \frac{m}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^m (r_i - \mu)^\top \Sigma^{-1} (r_i - \mu) \quad (*)$$

want to maximize observed log-likelihood:

$$\begin{aligned} l_{obs}(\theta) &= c - \frac{m_0}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^{m_0} (r_i - \mu)^\top \Sigma^{-1} (r_i - \mu) && \text{complete pairs of } x_i, y_i \\ &\quad - \frac{m_1}{2} \log \sigma_{22} - \frac{1}{2\sigma_{22}} \sum_{i=m_0+1}^{m_0+m_1} (y_i - \mu_2)^2 && \text{only } y_i \text{ observed} \\ &\quad - \frac{m_2}{2} \log \sigma_{11} - \frac{1}{2\sigma_{11}} \sum_{i=m_0+m_1+1}^m (x_i - \mu_1)^2 && \text{only } x_i \text{ observed} \end{aligned}$$

\rightarrow no closed solution to $\max l_{obs}(\theta)$

EM-algorithm: i) E-step: $E_{\tilde{\theta}}[l_m^c(\theta) | \text{observed d.}] = Q(\theta, \tilde{\theta})$

$\tilde{\theta}$ given depends only on $E_{\tilde{\theta}}[x_i | y_i], E_{\tilde{\theta}}[y_i | x_i]$ from (*)
 θ , and obs. data. $E_{\tilde{\theta}}[x_i^2 | y_i], E_{\tilde{\theta}}[y_i^2 | x_i]$

or by the fact that we have an exponential family of distributions.

$$x_i | y_i \sim N\left(\mu_1 + \frac{\sigma_{12}}{\sigma_{22}}(y_i - \mu_2); \sigma_{11} - \frac{\sigma_{12}\sigma_{22}}{\sigma_{22}}\sigma_{21}\right)$$

$$E_{\tilde{\theta}}[X_i | Y_i] = \tilde{\mu}_1 + \tilde{\sigma}_{12} \tilde{\sigma}_{22}^{-1} (Y_i - \tilde{\mu}_2)$$

$$E_{\tilde{\theta}}[X_i^2 | Y_i] = m_{\tilde{\theta}}(X_i | Y_i) + [E_{\tilde{\theta}}(X_i | Y_i)]^2$$

$$= \tilde{\sigma}_{11} - \tilde{\sigma}_{12} \tilde{\sigma}_{22}^{-1} \tilde{\sigma}_{21} + [E_{\tilde{\theta}}(X_i | Y_i)]^2$$

$$\tilde{\sigma}_{11} \left(1 - \frac{\tilde{\sigma}_{12}}{\tilde{\sigma}_{22} \tilde{\sigma}_{12}}\right) = \tilde{\sigma}_{11} (1 - \tilde{\rho}^2)$$

for $Z(Y_i | X_i)$ analogous.

ii) M-step: maximize $Q(\theta, \tilde{\theta})$ over θ

- for μ this leads to the usual MLE with replaced values of model obs.

$$\tilde{z}_i := (\tilde{x}_i, \tilde{y}_i) \text{ for } \tilde{x}_i = \begin{cases} x_i & \text{if } x_i \text{ is observed} \\ E_{\tilde{\theta}}[X_i | Y_i] & \text{if } x_i \text{ is not observed and } y_i \text{ anal.} \end{cases}$$

$$\tilde{\mu}_{\text{new}} = \frac{1}{n} \sum_{i=1}^n \tilde{z}_i = \frac{1}{n} \sum_{i=1}^n E[\tilde{z}_i | \text{obs}]$$

$$\left[\frac{\partial}{\partial \mu} \left[-\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum (\tilde{z}_i - \mu)' \Sigma^{-1} (\tilde{z}_i - \mu) \right] \right] = \sum (\tilde{z}_i - \mu)' \Sigma^{-1} = 0$$

$$\hat{\mu} = \frac{1}{n} \sum \tilde{z}_i$$

$$\text{for } \Sigma: \quad \frac{\partial}{\partial \Sigma} \left[-\frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum (\tilde{z}_i - \hat{\mu})' \underbrace{\Sigma^{-1} (\tilde{z}_i - \hat{\mu})}_{\text{Tr}((\tilde{z}_i - \hat{\mu})' \Sigma^{-1} (\tilde{z}_i - \hat{\mu}))} \right]$$

$$\text{Tr}((\tilde{z}_i - \hat{\mu})' \Sigma^{-1} (\tilde{z}_i - \hat{\mu})) = \text{Tr}(\Sigma' (\tilde{z}_i - \hat{\mu})(\tilde{z}_i - \hat{\mu})')$$

$$\left. \begin{array}{l} \frac{\partial}{\partial A} \text{Tr}(A) = I \\ \frac{\partial \text{Tr}(A)}{\partial A} = I A A^{-1} \\ \frac{\partial |A|}{\partial A} = |A| A^{-1} \end{array} \right\} \Rightarrow = -\frac{n}{2} \frac{|\Sigma| \Sigma^{-1}}{|\Sigma|} + \frac{1}{2} \sum \text{Tr}(\Sigma' (\tilde{z}_i - \hat{\mu})(\tilde{z}_i - \hat{\mu})') \Sigma' = 0$$

$$\Sigma = \frac{1}{n} \sum (\tilde{z}_i - \hat{\mu})(\tilde{z}_i - \hat{\mu})'$$

$$\frac{\partial \text{Tr}(A' B)}{\partial A} = -A' B A^{-1}$$

$$\left| \begin{array}{l} \tilde{\Sigma}_{\text{new}} = \frac{1}{n} \left[E[(z_i - \hat{\mu}_{\text{new}})(z_i - \hat{\mu}_{\text{new}})'] \mid \text{obs} \right] \\ = \frac{1}{n} \sum z_i z_i' \\ = \frac{1}{n} \sum E[z_i z_i' \mid \text{obs}] - \hat{\mu}_{\text{new}} \hat{\mu}_{\text{new}}' \\ = \frac{1}{n} \sum E[(x_{ij}^2, x_{ij} y_{ij}) \mid \text{obs}] - \hat{\mu}_{\text{new}} \hat{\mu}_{\text{new}}' \end{array} \right.$$

+ iterate EM-steps until convergence

Renweighting: if $P(X_i \text{ is not observed}) = r_i$ can be estimated, assign to each observed X_i a weight proportional to $\frac{1}{r_i}$ and perform analysis.