

Recall Example 5.2 (pp. 72–73)

X_1, X_2, \dots, X_n independent, each $\sim \text{Poisson}(\theta v_i)$

$$L(\theta) = \prod_{i=1}^n \frac{e^{-\theta v_i} (\theta v_i)^{x_i}}{x_i!} \quad \Rightarrow \quad \ell(\theta) = \sum_{i=1}^n -\theta v_i + x_i \log(\theta v_i) - \log x_i!$$

$$\ell'(\theta) = 0 \quad \Rightarrow \quad \hat{\theta} = \frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n v_i}$$

interesting “twist” (Example 5.5, p. 79–80)

X_1, X_2, \dots, X_n not fully observable

instead, observe only whether $X_i > 0$ or $X_i = 0$

MLE with Missing Data

Components

D_{obs} : observed data

D_{mis} : unobserved data (e.g., missing, latent, ...)

Θ : model parameters

MLE Want to

$$\max_{\Theta} \ell(\Theta; D_{obs}, D_{mis})$$

but “stuck” because D_{mis} not available.

EM Algorithm

E-Step Compute

$$Q(\Theta; D_{obs}, \hat{\Theta}_{t-1}) = \mathbb{E}_{D_{mis}|D_{obs}; \hat{\Theta}_{t-1}} [\ell(\Theta; D_{obs}, D_{mis})],$$

i.e., take best guess.

M-Step Compute

$$\hat{\Theta}_t = \arg \max_{\Theta} Q(\Theta; D_{obs}, \hat{\Theta}_{t-1}),$$

i.e., update parameter estimate.

Remarks E = expectation. M = maximization. Will skip theory about why it works, convergence, etc.

Continuation: Example 5.6 (pp. 84–85)

$$\ell(\theta) = \sum_{i=1}^n -\theta v_i + X_i \log(\theta v_i)$$

$$\mathbb{E}_{\mathbf{D}_{mis} | \mathbf{D}_{obs}; \hat{\theta}_{t-1}}[\ell(\theta)] = \sum_{i=1}^n -\theta v_i + [\mathbb{E}_{\mathbf{D}_{mis} | \mathbf{D}_{obs}; \hat{\theta}_{t-1}}(X_i)] \log(\theta v_i)$$

$$\mathbb{E}_{\mathbf{D}_{mis} | \mathbf{D}_{obs}; \hat{\theta}_{t-1}}(X_i) = \begin{cases} \mathbb{E}(\overset{\mathbf{D}_{mis}}{\downarrow} \tilde{X}_i \mid \underbrace{X_i > 0}_{\mathbf{D}_{obs}}; \hat{\theta}_{t-1}), & X_i > 0; \\ 0, & X_i = 0 \end{cases}$$

Exercise 5.11 (p. 85) Show that, if $X \sim \text{Poisson}(\lambda)$, then $\mathbb{E}(X | X > 0; \lambda) = \lambda / (1 - e^{-\lambda})$.

Two-Stage Generating Mechanism

- $Z_i \in \{1, 2, \dots, K\} \sim \text{multinomial}(\mathbf{1}; \pi_1, \pi_2, \dots, \pi_K)$
- $X_i | (Z_i = k) \sim f(x; \theta_k)$
 - here (and often), same family f , different parameters θ_k
 - in principle, can also do f_k , but need to be careful
- but label $Z_i \in \{1, 2, \dots, K\}$ **unobserved**

Components

- parameters Θ

$$\Theta = \{\theta_1, \dots, \theta_K, \pi_1, \dots, \pi_K\}$$

- observed data D_{obs}

$$\mathbf{X} = \{X_1, \dots, X_n\} \quad \text{or} \quad \{x_1, \dots, x_n\}$$

- unobserved “data” D_{mis}

$$\mathbf{Z} = \{Z_1, \dots, Z_n\}$$

MLE with Missing Data

Likelihood (pretend z_i available)

$$L(\Theta; \mathbf{X}, \mathbf{Z}) = \prod_i f(x_i, \overset{\text{NA}}{\downarrow} z_i) = \prod_i \left\{ \prod_k \left[\underbrace{f(x_i; \theta_k)}_{f(x_i|z_i)} \overset{f(z_i)}{\downarrow} \pi_k \right]^{I(z_i=k)} \right\}$$

Log-Likelihood (remember z_i not available)

$$\ell(\Theta; \mathbf{X}, \mathbf{Z}) = \sum_i \sum_k [I(Z_i = k)] \log [f(x_i; \theta_k) \pi_k]$$

MLE Would like to maximize $\ell(\Theta; \mathbf{X}, \mathbf{Z})$ over Θ but \mathbf{Z} missing.

EM Algorithm

E-Step Take a guess — suffices to compute

$$\mathbb{E}_{\mathbf{Z}|\mathbf{X};\hat{\Theta}_{t-1}}[I(Z_i = k)]$$

as log-likelihood is linear in $I(Z_i = k)$.

M-Step Update parameter estimate

$$\begin{aligned}\hat{\Theta}_t &= \arg \max_{\Theta} \sum_i \sum_k \underbrace{\mathbb{E}_{\mathbf{Z}|\mathbf{X};\hat{\Theta}_{t-1}}[I(Z_i = k)]}_{w_{ik}} \log [f(x_i; \theta_k) \pi_k] \\ &= \arg \max_{\Theta} \sum_i \sum_k w_{ik} \log [f(x_i; \theta_k)] + \sum_i \sum_k w_{ik} \log(\pi_k).\end{aligned}$$

M-Step

- for the “ θ_k part”, just **weighted** maximum likelihood

$$\max_{\theta_k} \sum_i w_{ik} \log [f(x_i; \theta_k)] \quad \forall k \quad (1)$$

- for the “ π_k part”, usual **multinomial** likelihood with w_{ik}

$$\max_{\pi_1, \dots, \pi_K} \sum_i \sum_k w_{ik} \log(\pi_k) \quad \text{s.t.} \quad \sum_k \pi_k = 1 \quad (2)$$

Remark A “prototypical” model here (for $\mathbf{x}_i \in \mathbb{R}^d$) is to take $f(\mathbf{x}_i; \boldsymbol{\theta}_k)$ as $N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$. [Fraley & Raftery (2002; *JASA*) discuss why the temptation $N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ is “problematic” and what can be done to compromise.]

E-Step

Simply writing $\hat{\Theta}$ rather than $\hat{\Theta}_{t-1}$,

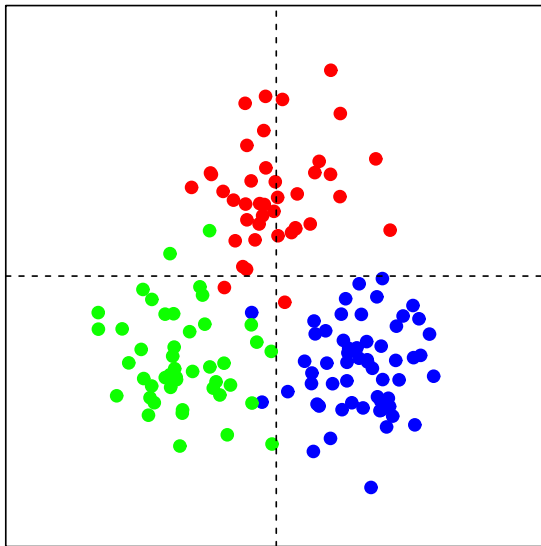
$$\begin{aligned}w_{ik} &\equiv \mathbb{E}_{\mathbf{Z}|\mathbf{X};\hat{\Theta}}[I(Z_i = k)] \\ &= \mathbb{P}(Z_i = k|X_i = x_i; \hat{\Theta}) \\ &= \frac{\hat{\pi}_k f(x_i; \hat{\theta}_k)}{\hat{\pi}_1 f(x_i; \hat{\theta}_1) + \dots + \hat{\pi}_K f(x_i; \hat{\theta}_K)}.\end{aligned}$$

The denominator shows the marginal distribution of each X_i is a [finite mixture distribution](#)

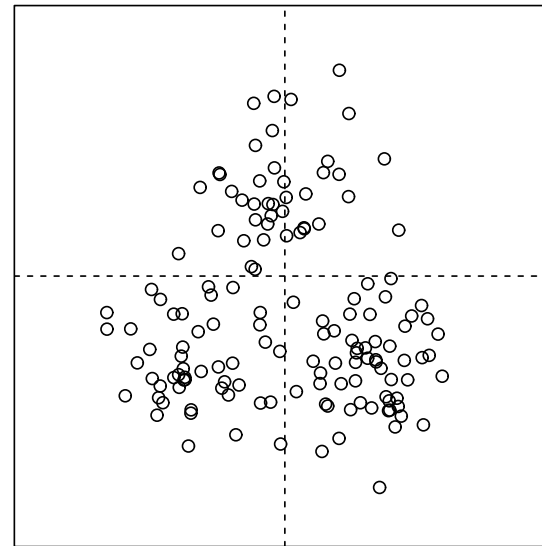
$$f(x_i; \Theta) = \sum_{k=1}^K \pi_k f(x_i; \theta_k).$$

Finite Mixture Model

(a)



(b)



(a) hidden generating mechanism; (b) observed data

Variation for Text Data

- each $\mathbf{x}_i = (x_{i1}, \dots, x_{id})^\top$ a text document
- x_{ij} = count of word j in document i (**bag of words**)
- can use the same framework, except choose

$$f(\mathbf{x}_i; \boldsymbol{\theta}_k) \text{ as } \text{multinomial}(m_i; \mathbf{p}_k),$$

where $\mathbf{p}_k = (p_{1k}, \dots, p_{dk})^\top$, and m_i = total number of words in document i (treated as fixed constants)

- many sophisticated variations/extensions, e.g., **latent Dirichlet allocation** (Blei, Ng & Jordan, **2003**; *JMLR*), mostly because vanilla MLE does not work well — most \hat{p}_{jk} would be zero

Back to Basketball Passes

- recall: given $\mathbf{Z} = \{z_1, \dots, z_n\}$,

$$\mathcal{M}(\mathbf{T}|\mathbf{Z}) = \prod_{i,j} \left\{ J(\rho_{z_i z_j}) \times \prod_{h=1}^{m_{ij}} \rho_{z_i z_j}(t_{ijh}) \right\}$$

- likelihood with latent variables Z_i, Z_j (not observed)

$$\prod_{i,j} \left\{ \prod_{k,\ell} \left[J(\rho_{k\ell}) \times \prod_{h=1}^{m_{ij}} \rho_{k\ell}(t_{ijh}) \right]^{I(Z_i=k, Z_j=\ell)} \right\} \times \prod_i \prod_k \pi_k^{I(Z_i=k)}$$

- requires $\mathbb{P}(Z_i = k, Z_j = \ell|\dots)$ in addition to $\mathbb{P}(Z_i = k|\dots)$, but Z_i, Z_j conditionally dependent here