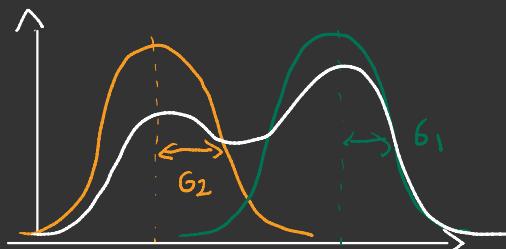


Modèles de Mélange fini

1894 . Karl Pearson .



$$f(x) = \pi_1 \underline{f_1(x)} + \pi_2 \underline{f_2(x)}$$

$$\frac{\mu_2}{\pi_2} < \frac{\mu_1}{\pi_1}$$

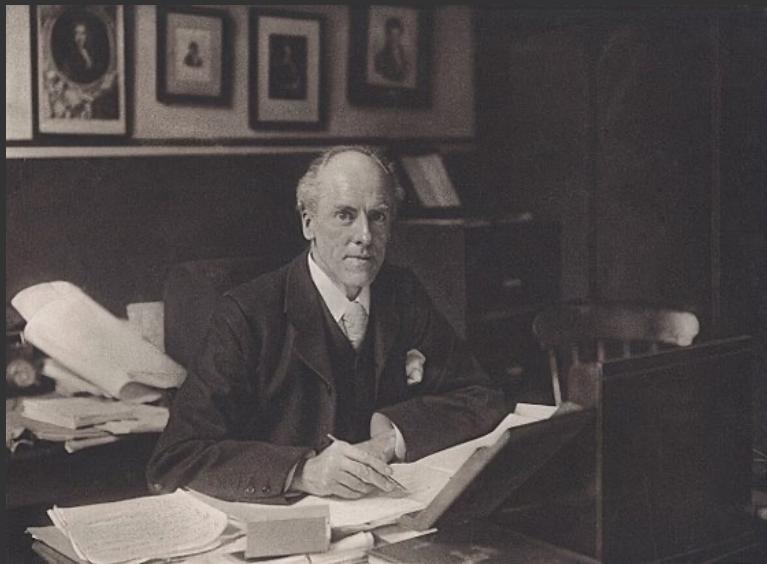
$$\forall x \quad f_m(x) = \frac{1}{\sqrt{2\pi}\sigma_p} \exp^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma_p}\right)^2}$$

$$\Theta = \left\{ \pi_1, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2 \right\}$$

$$\pi_2 = 1 - \pi_1$$

$$\sqrt{\pi_1} \in [0, 1]$$

KARL PEARSON 1894 (Fisher 1912)



KARL PEARSON

$$f(\alpha) = \sum_{h=1}^k \pi_h f_h(\alpha)$$

où les f_h
sont des densités
paramétriques.

III. "Contributions to the Mathematical Theory of Evolution."
By KARL PEARSON, M.A., Professor of Applied Mathematics,
University College. Communicated by Professor HENRICI,
F.R.S. Received October 18, 1893.

(Abstract.)

1. If a series of measurements, physical, biological, anthropological, or economical, not of the same object, but of a group of objects of the same type or family, be made, and a curve be constructed by plotting up the number of times the measurements fall within a given small unit of range to the range, this curve may be termed a *frequency curve*. As a rule this frequency curve takes the well known form of the curve of errors, and such a curve may be termed a *normal frequency curve*. The latter curve is symmetrical about its maximum ordinate. Occasionally, however, frequency curves do not take the normal form, and are then generally, but not necessarily, asymmetrical. Such abnormal curves arise particularly in biological measurements; they have been found by Professor

$$f(x) = \pi_1 f_1(x) + \pi_2 f_2(x)$$

Observées : $(\alpha_1, \dots, \alpha_n)$

Missingues : (y_1, \dots, y_2)

INITIALISATION :

$$\Theta^{(0)} = \left\{ \pi_1^{(0)}, \mu_1^{(0)}, G_1^{(0)}, \mu_2^{(0)}, G_2^{(0)} \right\}$$

$$\forall k \quad \pi_k^{(0)} = \frac{1}{K} \quad \text{ici} \quad \pi_1 = \frac{1}{2}$$

$$\forall k \quad \mu_k^{(0)} = \alpha_i \quad \text{ou } i \sim U\{1, \dots, n\}$$

$$\forall m \quad G_m = 1$$

ITERATION :

$$E : Q(\theta, \theta^q) = \mathbb{E}_{Z|X} \left[\log P(X, Z; \theta) \right]$$

$$M : \theta^{q+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta^q)$$

$$\begin{aligned} P(\alpha_1, \dots, \alpha_n, z_1, \dots, z_n; \theta) &= \prod_{i=1}^n P(\alpha_i, z_i; \theta) \\ &= \prod_{i=1}^n \left(\prod_{k=1}^K P(\alpha_i, z_i) \right)^{\Delta\{z_i = b_k\}} \end{aligned}$$

$$\log P(X, Z; \theta) = \sum_{i=1}^n \sum_{k=1}^K \underbrace{\mathbb{E}_{\{Z_i=k\}}}_{\log P(Z_i=k; \theta)} \underbrace{\log P(x_i | z_i=k)}_{\pi_k f_k(x_i; \theta)}$$

$$\mathbb{E}_{Z|X; \theta^q} \left[\log P(X, Z; \theta) \right] = \sum_i \sum_k \underbrace{\mathbb{E} \left[\mathbb{E}_{\{Z_i=k\}} \right]}_{P(Z_i=k|x_i; \theta^q)} \log \pi_k f_k(x_i; \theta)$$

$$T = \begin{bmatrix} & \dots & K \\ \vdots & \left[\begin{array}{c} \hline \\ \hline \end{array} \right] \\ \vdots & \end{bmatrix} \quad t_{ik}^{(q)} = \frac{1}{\sqrt{2\pi} G_k} \exp \left(-\frac{1}{2} \frac{(x_i - \mu_k)^2}{G_k^2} \right)$$

$$M_1 \frac{\partial Q(\theta, \theta^q)}{\partial \mu_k} = 0$$

$$0 = \frac{\partial}{\partial \mu_k} \left(\sum_i \sum_k t_{ik}^{(q)} \left(\log \pi_k + \log \left(\frac{1}{\sqrt{2\pi} G_k} \exp \left(-\frac{1}{2} \frac{(x_i - \mu_k)^2}{G_k^2} \right) \right) \right) \right)$$

$$0 = \frac{\partial \sum_{i=1}^n t_{i,h}^{(q)} \left[\frac{1}{2} \left(\frac{x_i - \mu_h}{\sigma_h^2} \right)^2 \right]}{\partial \mu_h} \Rightarrow \sum_{i=1}^n t_{i,h}^{(q)} (x_i - \mu_h) \times 1 = 0$$

$$\partial \mu_h$$

$$\sum_{i=1}^n t_{i,h}^{(q)} x_i = \sum t_{i,q}^{(q)} \mu_h$$

$$\mu_h^{q+1} = \frac{\sum_{i=1}^n t_{i,h}^{(q)} x_i}{\sum_{i=1}^n t_{i,h}^{(q)}}$$

$$\frac{\partial Q(\theta, \theta^q)}{\partial \theta^q} = 0 \Rightarrow$$

$$\frac{\partial \left(\sum_{i=1}^n t_{i,h}^{(q)} \log \frac{1}{\sqrt{2\pi} \sigma_h^2} \exp^{-\frac{1}{2} \left(\frac{(x_i - \mu_h)^2}{\sigma_h^2} \right)} \right)}{\partial \sigma_h^2} = 0$$

$$\partial \sigma_h^2$$

$$\frac{\partial \left(\sum_i t_{i,h}^{(q)} \left(-\frac{1}{2} \log \sigma_h^2 \right) \right) + \partial \left(\sum_i t_{i,h}^{(q)} \left(-\frac{1}{2} \left(\frac{(x_i - \mu_h)^2}{\sigma_h^2} \right) \right) }{\partial \sigma_h^2} = 0$$

$$\partial \sigma_h^2$$

$$-\cancel{\frac{1}{2} \sum_i t_{i,h}^{(q)} \times \frac{1}{\sigma_h^2}} - \cancel{\frac{1}{2} \sum_i t_{i,h}^{(q)} \left(-\frac{1}{(\sigma_h^2)^2} (x_i - \mu_h)^2 \right)} = 0$$

$$G_{\theta}^{q+1} = \frac{\sum_{i=1}^n t_{i,h}^q \left(\theta_i - \pi_h^{q+1} \right)^2}{\sum_{i=1}^n t_{i,h}^q}$$

$$\sum_{h=1}^K \pi_h = 1 \quad \pi_1 + \pi_2 = 1$$

$$\mathcal{L}(\theta, \lambda) = Q(\theta, \theta^q) + \lambda \left(\sum_h \pi_h - 1 \right)$$

$$\begin{cases} \frac{\partial \mathcal{L}(\theta, \lambda)}{\partial \theta} = 0 & \left\{ \frac{\partial \left(\sum_{i=1}^n t_{i,h}^q \log \pi_h \right) + \partial (\lambda \sum_h \pi_h)}{\partial \pi_h} = 0 \right. \\ \left. \sum \pi_h - 1 = 0 \right. \end{cases}$$

$$\begin{cases} \sum_{i=1}^n t_{i,h}^{(q)} \times \frac{1}{\pi_h} - \lambda = 0 \\ \sum \pi_h - 1 = 0 \end{cases} \quad \left\{ \pi_h^{(q+1)} = \frac{\sum_{i=1}^n t_{i,h}^q}{n} \right.$$

$$\lambda = \sum_{i=1}^n \left(\sum_{h=1}^K t_{i,h}^q \right) = \sum_{i=1}^n 1 = n$$

$$\boxed{\pi_h^{q+1} = \frac{\sum_{i=1}^n t_{i,h}^q}{n}}$$

CONSEILS DE PROGRAMMATION

$$P(X=\alpha) = \sum_{k=1}^K \pi_k P_k(X=\alpha)$$

$$P_k(X=\alpha) = \frac{e^{-\lambda} \lambda^\alpha}{\alpha!} \quad \mathbb{E}[X | Z=k] = \lambda_k$$

$$\text{Init}(\underbrace{x_1, \dots, x_n}_X, K) \longrightarrow \Theta^{(0)} = \left\{ \pi_1^{(0)}, \pi_{K+1}^{(0)}, \lambda_1^{(0)}, \lambda_K^{(0)} \right\}$$

$$\begin{cases} E\text{-step}(\Theta^{(q)}, x_1, \dots, x_n) \longrightarrow T^{(q)} = (t_{ik})_{\substack{i=1 \dots n \\ k=1 \dots K}} \\ M\text{-step}(T, x_1, \dots, x_n) \longrightarrow \Theta^{(q+1)} \end{cases}$$

CRITÈRE DE CONVERGENCE

$$\frac{\|\Theta^q - \Theta^{q+1}\|_2^2}{\|\Theta^q\|_2^2} < \epsilon \quad (\text{ex } \epsilon = 10^{-6})$$

ALGORITHME EM POUR LES MODÈLES DE MÉLANGE FINI

Dempster, Laird, Rubin (1977)

OBSERVATIONS $(\alpha_1, \dots, \alpha_n)$

MANQUANTES (z_1, \dots, z_n)

COMPLÈTES $((\alpha_1, z_1), \dots, (\alpha_n, z_n))$

$$f(\alpha) = \sum_{h=1}^k \pi_h f_h(\alpha)$$

\uparrow proportion \uparrow densité de composante

$$\ell(\theta; \underbrace{\alpha_1, \dots, \alpha_n}_x) = \prod_{i=1}^n f(\alpha_i) = \prod_{i=1}^n \left(\sum_{h=1}^k \pi_h f_h(\alpha_i) \right)$$

$$\log \ell(\theta; \underbrace{\alpha_1, \dots, \alpha_n}_x) = \sum_{i=1}^n \log \left(\sum_{h=1}^k \pi_h f_h(\alpha_i) \right)$$

$$P(\underbrace{(\alpha_1, z_1), \dots, (\alpha_n, z_n)}_X) = P(X, Z) = P(X) \times P(Z|X)$$

$$\log \underbrace{\ell(\theta; X)}_{P(X)} = \log P(X, Z; \theta) - \log P(Z|X; \theta)$$

$$\underbrace{\mathbb{E}_{Z|X;\theta^q} \left[\log P(X; \theta) \right]}_{\text{Q}(\theta, \theta^q)} = \boxed{\mathbb{E}_{Z|X;\theta^q} \left[\log P(X, Z; \theta) \right]} - \text{H}(\theta, \theta^q)$$

$$\boxed{\begin{aligned} & \log P(X; \theta) \\ &= \\ & Q(\theta, \theta^q) - H(\theta, \theta^q) \end{aligned}}$$

$$\boxed{\mathbb{E}_{Z|X;\theta^q} \left[\log P(Z|X; \theta) \right]}$$

INIT : $\theta^{(0)}$

REPETER JUSQU'à conv.

$E(\text{expectation})$: CALCUL DE $Q(\theta, \theta^q)$
$M(\text{maximisation})$: $\theta^{q+1} = \underset{\theta}{\operatorname{argmax}} Q(\theta, \theta^q)$

MONTRONS QUE CHAQUE ITERATION DE EN
AUGMENTE $\ell(\theta; x)$

$$Q(\theta^{q+1}, \theta^q) \geq Q(\theta^q, \theta^q)$$

$$\times \mathbb{E}_{Z|X, \theta^q} \left[\log P(X, Z; \theta^{q+1}) - \log P(X, Z; \theta^q) \right] \geq 0$$

$$\int P(Z|X; \theta^q) \log \frac{P(X, Z; \theta^{q+1})}{P(X, Z; \theta^q)} dZ \geq 0$$

$$\log \mathbb{E}_{Z|X; \theta^q} \left[\frac{P(X, Z; \theta^{q+1})}{P(X, Z; \theta^q)} \right] \geq 0$$

SENSEN

$$\log \int P(Z|X; \theta^q) \frac{P(X, Z; \theta^{q+1})}{P(X, Z; \theta^q)} dZ \geq 0$$

$$\log \int P(Z|X; \theta^q) \frac{P(X; \theta^{q+1})}{P(X; \theta^q)} \frac{P(Z|X; \theta^{q+1})}{P(Z|X; \theta^q)} dZ \geq 0$$

$$\log \frac{P(X; \theta^{q+1})}{P(X; \theta^q)} \int P(Z|X; \theta^{q+1}) dZ \geq 0 \Leftrightarrow \boxed{\frac{P(X; \theta^{q+1})}{P(X; \theta^q)} \geq 1}$$