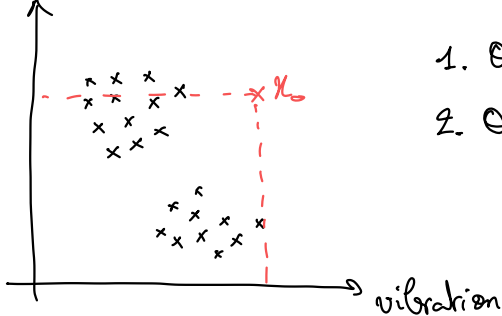


III) Algorithme EM appliqué aux mélanges gaussiens

Exemple d'application de l'estimation de modèle

température



1. On estime un modèle $p(x)$

2. On met en place une règle de décision
si $p(x_0) < \epsilon \rightarrow$ anomalie

Exemple introduit "simple"

Estimation par MIV d'une gaussienne simple ($K=1$)

$$X = (x_1, \dots, x_N)^T \text{ iid } x_n \sim \mathcal{N}(\mu^*, \Sigma^*) \quad \forall n \in \llbracket 1, N \rrbracket$$

On cherche $\theta_{ML} = (\mu_{ML}, \Sigma_{ML})$ qui maximisent la
log-vraisemblance $\ell(\theta) = \ln(p(X|\theta)) = \ln\left(\prod_{n=1}^N p(x_n|\theta)\right)$
 $= \sum_{n=1}^N \ln(p(x_n|\theta)).$

En utilisant la pdf d'une gaussienne multivariée ($D>1$):

$$\mathcal{N}(x|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

On obtient :

$$\underbrace{\ell(\mu, \Sigma)}_{=\theta} = \frac{-ND}{2} \ln(2\pi) - \frac{N}{2} \ln(|\Sigma|) - \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^T \Sigma^{-1} (x_n - \mu)$$

Lemmes : $\det(A)$

$$\textcircled{A} \quad \frac{\partial \ln(|A|)}{\partial A_{ij}} = (A^{-1})_{ji}$$

$$\textcircled{B} \frac{\partial A^{-1}}{\partial A_{ij}} = -A^{-1} E_{ij} A^{-1} \quad (E_{ij} = \delta_{ij}) \quad i \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}^{-1} \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\textcircled{C} \text{Tr}(A E_{ij}) = A_{ji}$$

$$\textcircled{D} g: x \mapsto (x | Ax), \quad \frac{\partial f}{\partial x} = (A + A^T)x.$$

Dérivation de $\ell(\mu, \Sigma)$:

1. Estimation de μ_{MLE} :

$$\frac{\partial \ell}{\partial \mu} \textcircled{D} = -\frac{1}{2} \sum_{n=1}^N (-2) \Sigma^{-1} (x_n - \mu)$$

$$= \Sigma^{-1} \left(\sum_{n=1}^N (x_n - \mu) \right)$$

$$\text{donc } \frac{\partial \ell}{\partial \mu} = 0 \Leftrightarrow \Sigma^{-1} \left(\sum_{n=1}^N (x_n - \mu) \right) = 0$$

$$\stackrel{\Sigma^{-1} \text{ inv.}}{\Leftrightarrow} \left(\sum_{n=1}^N x_n \right) - N\mu = 0$$

$$\Leftrightarrow \mu = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\text{Donc } \boxed{\mu_{MLE} = \frac{1}{N} \sum_{n=1}^N x_n}$$

2. Estimation de Σ_{MLE}

$$\frac{\partial \ell}{\partial \Sigma_{ij}} \textcircled{A, B} = -\frac{N}{2} (\Sigma^{-1})_{ji} + \frac{1}{2} \sum_{n=1}^N \overbrace{(x_n - \mu)^T \Sigma^{-1} E_{ij} \Sigma^{-1} (x_n - \mu)}^{(1, D) \quad (D, 1)}$$

$$:= S_n \quad (1, 1)$$

On a $S_n \in \mathbb{R}$ donc $S_n = \text{Tr}(S_n)$,

$$\text{ie } (x_n - \mu)^T \Sigma^{-1} E_{ij} \Sigma^{-1} (x_n - \mu) = \text{Tr}((x_n - \mu)^T \Sigma^{-1} E_{ij} \Sigma^{-1} (x_n - \mu))$$

$$= \text{Tr}(\Sigma^{-1} (x_n - \mu) (x_n - \mu)^T \Sigma^{-1} E_{ij})$$

$$= (\Sigma^{-1} (x_n - \mu)(x_n - \mu)^T \Sigma^{-1})_{ji}$$

$$\text{Donc } \frac{\partial \ell}{\partial \Sigma_{ij}} = \left[-\frac{N}{2} \Sigma^{-1} + \frac{1}{2} \sum_{n=1}^N \Sigma^{-1} (x_n - \mu)(x_n - \mu)^T \Sigma^{-1} \right]_{ji}$$

$$\text{Enfin, } \left[\frac{\partial \ell}{\partial \Sigma_{ij}} = 0 \quad \forall i, j \right] \Leftrightarrow -\frac{N}{2} \Sigma^{-1} + \frac{1}{2} \sum_{n=1}^N \Sigma^{-1} (x_n - \mu)(x_n - \mu)^T \Sigma^{-1} = 0$$

multiplier à g. et à d. par Σ .

$$\Leftrightarrow -\frac{N}{2} \Sigma + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T = 0$$

$$\Leftrightarrow \Sigma = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T$$

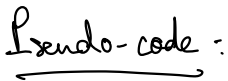
$$\text{Donc } \boxed{\Sigma_{ML} = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T}$$

Estimation pour un mélange de K gaussiennes

Contexte: On suppose que des observations $x_1, \dots, x_N \in \mathbb{R}^D$ sont issues d'un mélange gaussien de paramètres $\theta^* = (\pi^*, \mu^*, \Sigma^*)$.

On cherche θ qui maximise la log-vraisemblance

$$\ell(\theta) = \ln \left(\prod_{n=1}^N p(x_n | \theta) \right)$$



Répétier :

- E-step: calculer les postérieurs $\gamma_{nk} = p(z_n = k | x = x_n)$ avec les valeurs actuelles de θ pour en déduire la fonction $Q(\theta)$.
- M-step: mise à jour des paramètres $\theta = (\pi, \mu, \Sigma)$ pour maximiser $Q(\theta)$.

Critère d'arrêt: par exemple quand les paramètres ont "convergé".

On tente un maximum vraisemblance pour un mélange de K gaussiennes

$X = (x_1, \dots, x_N)^T$: observations ($x_n \in \mathbb{R}^D \forall n \in \llbracket 1, N \rrbracket$)

$Z = (z_1, \dots, z_N)$: variables ($z_n \in \llbracket 1, K \rrbracket \forall n \in \llbracket 1, N \rrbracket$)

$$\theta = \begin{cases} \pi = (\pi_1, \dots, \pi_K) \text{ coefficients de m\u00e9lange } \pi_k = p(z=k) \in]0,1[\\ \text{et } \sum_{k=1}^K \pi_k = 1 \\ \mu = (\mu_1, \dots, \mu_K) \text{ moyennes des } K \text{ gaussiennes } (\mu_k \in \mathbb{R}^D) \\ \Sigma = (\Sigma_1, \dots, \Sigma_K) \text{ matrices de covariance des } K \text{ gaussiennes } \\ (\Sigma_k \in \mathbb{R}^{D \times D}) \end{cases}$$

On cherche les valeurs de π, μ, Σ qui maximisent :

$$\ell(\pi, \mu, \Sigma) = \ln \left(\prod_{n=1}^N p(x_n | \pi, \mu, \Sigma) \right)$$

$$\begin{aligned} &= \sum_{n=1}^N \ln(p(x_n | \pi, \mu, \Sigma)) \\ &\stackrel{\text{formule des proba totales}}{=} \sum_{n=1}^N \ln \left(\sum_{k=1}^K p(z_n=k, x_n | \pi, \mu, \Sigma) \right) \\ &= \sum_{n=1}^N \ln \left(\sum_{k=1}^K \underbrace{p(z_n=k)}_{\pi_k} \underbrace{p(x_n | z_n=k, \mu_k, \Sigma_k)}_{\mathcal{N}(x_n | \mu_k, \Sigma_k)} \right) \\ &= \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \underbrace{\mathcal{N}(x_n | \mu_k, \Sigma_k)}_{\frac{1}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2}(x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right)} \right) \end{aligned}$$

Probl\u00e8me : On ne peut pas diff\u00e9rencier par rapport aux μ_k et Σ_k (log de somme d'exponentielle).

Idee: Rendre la log-vraisemblance séparable.

Rappel: Si $f: (x, y) \mapsto g(x) + h(y)$
alors $\arg \max_{(x, y)} f = (\arg \max_x g, \arg \max_y h)$

$$\ell(\theta) = \ln(p(X|\theta))$$

On introduit la log-vraisemblance "complétée":

$$\begin{aligned}\ell'(\theta) &= \ln(p(X, Z|\theta)) = \sum_{n=1}^N \ln(p(x_n, z_n|\theta)) \\ &= \sum_{n=1}^N \ln(p(z=z_n) p(x=x_n | z=z_n, \theta)) \\ &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{z_n=k} \ln(\underbrace{p(z=k)}_{=\pi_k} \underbrace{p(x=x_n | z=k, \theta)}_{=\mathcal{N}(x_n | \mu_k, \Sigma_k)}) \\ &= \sum_{n=1}^N \sum_{k=1}^K \mathbb{1}_{z_n=k} \ln(\pi_k \underbrace{\mathcal{N}(x_n | \mu_k, \Sigma_k)}_{\substack{\frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp(\dots)}})\end{aligned}$$

Pour maximiser cette quantité $\ell'(\theta)$, on va maximiser son espérance (heuristique \S):

$$\mathbb{E}(\ell'(\theta)) = \sum_{n=1}^N \sum_{k=1}^K \mathbb{E} \left[\mathbb{1}_{z_n=k} \ln(\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)) \right]$$

Lemme: Si X, Y 2 v.a.r. $\mathcal{A}_q(X, Y)$ admet une espérance sur \mathbb{R}^2 , alors $\mathbb{E}_X(X) = \mathbb{E}_Y(\mathbb{E}_{X|Y}(X|Y))$.

Preuve: $\mathbb{E}_Y(\mathbb{E}_X(X|Y)) = \int_Y \mathbb{E}_X(X|Y=y) f_Y(y) dy$

$$= \int_y \left[\int_x x f_{x|y}(x) dx \right] f_y(y) dy$$

Fubini (cf. ci-dessous)

$$= \int_x x \left[\int_y \underbrace{f_{x|y}(x) f_y(y)}_{= f_{(x,y)}(x,y)} dy \right] dx$$

$$\Rightarrow [P(X|Y)P(Y) = P(X,Y)]$$

$$= \int_x x f_x(x) dx$$

$$= E_x(X)$$

$$E_{X,Y}(x,y) = \iint_{(u,v)} f_{x,y}(u,v) \frac{dx dy}{du dv}$$

Fubini: $\underbrace{|x f_{x|y}(x) f_y(y)|}_{= f_{(x,y)}(x,y)} \leq \underbrace{\| \begin{pmatrix} x \\ y \end{pmatrix} \| \times |f_{(x,y)}(x,y)|}_{\text{intégrable car } (X,Y) \text{ admet une espérance}}$

CQFD

On a donc

$$E_{\left(\begin{smallmatrix} X \\ Y \end{smallmatrix} \right)}(l'(\theta)) = E_X \left(E_{\left(\begin{smallmatrix} Y \\ \theta \end{smallmatrix} \right)}(l'(\theta) | X) \right)$$

$$= \sum_{n=1}^N \sum_{k=1}^K E_{x_n} \left[E \left[\mathbb{1}_{z_n=k} \left(\ln(\pi_k') + \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k')) \right) \mid x=x_n \right] \right]$$

$$= \sum_{n=1}^N \sum_{k=1}^K E_{x_n} \left[\underbrace{E \left(\mathbb{1}_{z_n=k} \mid x=x_n \right)}_{= p(z_n=k | x=x_n)} \left[\ln(\pi_k') + \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k')) \right] \right]$$

$$= \sum_{n,k} \dots$$

$$= E_{\mathbf{X}} \left[\underbrace{\sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} [\ln(\pi_k) + \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k))]}_{:= Q(\theta)} \right]$$

Q est bien séparable à γ_{nk} fixés : en fait on va considérer que les postérieurs γ_{nk} (qui ont été calculés avec les valeurs actuelles des paramètres) sont fixés.

Retour sur les idées à l'origine de l'expression de Q :

Quand on a une seule gaussienne, on maximise $l(\theta) = \ln(p(\mathbf{X}|\theta))$. Pour un mélange à K gaussiennes,

$$l(\theta) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k) \right)$$

↳ pas tractable

Idee : utiliser la log-vraisemblance "complète" :

$$p(\mathbf{X}, \mathbf{Z} | \theta) = p(\mathbf{X} | \theta) p(\mathbf{Z} | \mathbf{X}, \theta)$$

$$\text{ie } \underbrace{\ln(p(\mathbf{X}, \mathbf{Z} | \theta))}_{\substack{\approx Q(\theta) \\ \text{moment}}} = \underbrace{\ln(p(\mathbf{X} | \theta))}_{l(\theta)} + \underbrace{\ln(p(\mathbf{Z} | \mathbf{X}, \theta))}_{\leq 0}$$

Pour formaliser, on introduit la distribution $q(\mathbf{Z})$:

$$l(\theta) = \ln(p(\mathbf{X} | \theta)) = \ln(p(\mathbf{X}, \mathbf{Z} | \theta)) - \ln(p(\mathbf{Z} | \mathbf{X}, \theta))$$

$$= \ln\left(\frac{p(x, z | \theta)}{q(z)}\right) - \ln\left(\frac{p(z | x, \theta)}{q(z)}\right)$$

→ $xq(z)$

→ \sum_z

↑

$$= \underbrace{\sum_z q(z) \ln\left(\frac{p(x, z | \theta)}{q(z)}\right)}_{\mathcal{L}(q, \theta)} + \underbrace{\sum_z -q(z) \ln\left(\frac{p(z | x, \theta)}{q(z)}\right)}_{= KL(q \| p(z | x, \theta))}$$

$$\sum_z q(z) \times \text{cst} = \text{cst} \left(\sum_z q(z) \right)$$

$\mathcal{L}(q, \theta)$

$= KL(q \| p(z | x, \theta))$

Rappel sur la Kullback-Leibler

Si P et Q sont 2 star de densités p et q ,

$$\text{alors } KL(P \| Q) = \int_{-\infty}^{+\infty} p(x) \ln\left(\frac{p(x)}{q(x)}\right) dx = E_P\left(\ln\left(\frac{p}{q}\right)\right)$$

"KL de P par Q à Q "

Rmq: • mesure de l'information qu'on perdrait si on utilisait Q pour approximer P

• KL n'est pas symétrique

• $KL(P \| Q) \geq 0$ avec \Leftrightarrow si $P = Q$

Donc on réécrit ce qu'on avait :

$$\ell(\theta) = \mathcal{L}(q, \theta) + KL(q \| p(z | x, \theta))$$

E-step: On cherche à maximiser $\mathcal{L}(q, \theta^{\text{old}})$ en fct° de q ,

or $\ell(\theta)$ ne dépend pas de q , donc cela revient à minimiser $KL(q \| p(z | x, \theta))$, ie de forcer $q = p(z | x, \theta^{\text{old}})$.

M-step: $\ell(\theta) = \mathcal{L}(p(z | x, \theta^{\text{old}}), \theta)$

$$= \sum_z p(z|x, \theta^{old}) \ln \left(\frac{p(x, z | \theta)}{p(z|x, \theta^{old})} \right)$$

$$= \underbrace{\sum_z p(z|x, \theta^{old}) \ln(p(x, z | \theta))}_{Q(\theta)} - \underbrace{\sum_z p(z|x, \theta^{old}) \ln(p(z|x, \theta^{old}))}_{\text{ne dépend pas de } \theta}$$

On cherche à maximiser :

$$Q(\theta) = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \left[\ln(\pi_k) + \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k)) \right]$$

M-step en calculs :

① Maximisation de Q en fct^o de π

On doit respecter $\sum_{k=1}^K \pi_k = 1$ donc il s'agit d'une optimisation sous contrainte, donc on introduit le lagrangien

$$\mathcal{L}(\pi, \lambda) = Q(\pi, \mu, \Sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \sum_{n=1}^N \frac{\gamma_{nk}}{\pi_k} + \lambda = \frac{1}{\pi_k} \sum_{n=1}^N \gamma_{nk} + \lambda$$

$$\left[\forall k \in [1, K], \frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \right] \Rightarrow \sum_{k=1}^K \pi_k \frac{\partial \mathcal{L}}{\partial \pi_k} = 0$$

$$\Rightarrow \sum_{k=1}^K \sum_{n=1}^N \gamma_{nk} + \underbrace{\sum_{k=1}^K \pi_k}_{=1} \lambda = 0$$

$$\Rightarrow \sum_{n=1}^N \underbrace{\sum_{k=1}^K \gamma_{nk}}_{p(z_n=k | x_n)} + \lambda = 0$$

$$= 1$$

$$\Rightarrow \lambda = -N$$

En réinjectant la valeur de λ dans $\frac{\partial \mathcal{L}}{\partial \pi_k}$,
on obtient $\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{1}{\pi_k} \sum_{n=1}^N \gamma_{nk} - N$.

$$\text{Donc } \frac{\partial \mathcal{L}}{\partial \pi_k} = 0 \Leftrightarrow \pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}$$

$$\boxed{\pi_k^{\text{step}} = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}}$$

Réq: C'est la moyenne sur toutes les observations de leur probabilité (postérieure) d'appartenir à la $k^{\text{ème}}$ gaussienne.

② Maximisation de Q-en fct° de μ

$$\ln(\mathcal{N}(x_n, \mu_k, \Sigma_k)) = -\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln(|\Sigma_k|) - \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)$$

$$\frac{\partial \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k))}{\partial \mu_k} = \Sigma_k^{-1} (x_n - \mu_k)$$

$f: x \mapsto \langle x | A x \rangle$
 $\Rightarrow \nabla f_x = (A + A^T) x$
 et ici Σ_k^{-1} symétrique

$$\begin{aligned} \frac{\partial Q}{\partial \mu_k} &= \sum_{n=1}^N \gamma_{nk} \Sigma_k^{-1} (x_n - \mu_k) \\ &= \Sigma_k^{-1} \left(\sum_{n=1}^N \gamma_{nk} x_n - \mu_k \sum_{n=1}^N \gamma_{nk} \right) \end{aligned}$$

donc $\frac{\partial Q}{\partial \mu_k} = 0 \Leftrightarrow \boxed{\mu_k = \frac{\sum_{n=1}^N \gamma_{nk} x_n}{\sum_{n=1}^N \gamma_{nk}}}$

Remq: C'est la moyenne de toutes les observations pondérées par leur probabilité (postérieure) d'appartenir à la même gaussienne.

③ Maximisation en fct° de Σ

Lemmes: (a) $\frac{\partial \ln(|A|)}{\partial A_{ij}} = (A^{-1})_{ji}$

(b) $\frac{\partial A^{-1}}{\partial A_{ij}} = -A^{-1} E_{ij} A^{-1} \quad (E_{ij} = \delta_{ij})$

(c) $\text{Tr}(A E_{ij}) = A_{ji}$

$$\frac{\partial \ln(\mathcal{N}(x_n | \mu_k, \Sigma_k))}{\partial \Sigma_k} =$$

$$\frac{\partial \Sigma_k}{\partial \Sigma_k, ij} = \frac{1}{2} (\Sigma_k^{-1})_{ji} + \frac{1}{2} \underbrace{(x_n - \mu_k)^T \Sigma_k^{-1} E_{ij} \Sigma_k^{-1} (x_n - \mu_k)}_{1 \times 1}$$

$$= \frac{1}{2} (\Sigma_k^{-1})_{ji} + \frac{1}{2} \text{Tr} \left(\underbrace{(x_n - \mu_k)^T \Sigma_k^{-1} E_{ij} \Sigma_k^{-1} (x_n - \mu_k)}_{1 \times 1} \right)$$

$$= \frac{1}{2} (\Sigma_k^{-1})_{ji} + \frac{1}{2} \text{Tr} \left(\Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1} E_{ij} \right)$$

$$= \frac{1}{2} (\Sigma_k^{-1})_{ji} + \frac{1}{2} \left[\Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1} \right]_{ji}$$

$$= \frac{1}{2} \left[\Sigma_k^{-1} + \Sigma_k^{-1} (x_n - \mu_k) (x_n - \mu_k)^T \Sigma_k^{-1} \right]_{ji}$$

$$\Rightarrow \frac{\partial \ln(N(x_n | \mu_k, \Sigma_k))}{\partial \Sigma_k} = \frac{1}{2} \left[\Sigma_k^{-1} + \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} \right]$$

$$\frac{\partial Q}{\partial \Sigma_k} = 0 \quad (\Rightarrow) \quad \sum_{n=1}^N \gamma_{nk} \left[\Sigma_k^{-1} + \Sigma_k^{-1} (x_n - \mu_k)(x_n - \mu_k)^T \Sigma_k^{-1} \right] = 0$$

$$\Rightarrow \sum_{n=1}^N \gamma_{nk} \Sigma_k + \sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T = 0$$

$$\Rightarrow \boxed{\Sigma_k = \frac{\sum_{n=1}^N \gamma_{nk} (x_n - \mu_k)(x_n - \mu_k)^T}{\sum_{n=1}^N \gamma_{nk}}}$$