

L'algorithme EM et ses variantes:

On va se concentrer maintenant sur un autre algorithme itératif aussi qui possède de nombreux avantages : sa forme initiale est très utilisée et possède, sous des conditions réalistes, des propriétés de convergence. De plus, quand le problème se complexifie un peu, plusieurs variantes sont disponibles permettant de garder des propriétés de convergence intéressantes et d'avoir une implémentation facile et efficace.

Cet algorithme est l'algorithme EM pour Expectation - Maximisation. On va d'abord voir comment il est apparu puis ses propriétés et enfin ses variantes.

1 Rappel du cadre:

(Dempster, Laird, Rubin; 77)

On cherche à maximiser une quantité se mettant sous la forme d'une espérance :

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmax}} q(y; \theta) = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_z [q(y, z; \theta)] \quad (*)$$

où y est une donnée que l'on observe, typiquement un n -uplet de vecteurs (y_1, \dots, y_n) et z est une variable aléatoire de loi $q(z; \theta)$. Notez que cette loi peut dépendre de θ .

Rq: • Si $n=1$, on se retrouve dans le cas où on a à maximiser $\int q(y, z; \theta) d\mu(z)$ ce qui se ramène au cas considéré précédemment.

• Si $n \neq 1$ mais (y_i) iid, on se retrouve aussi avec $\int (\sum q(y_i, z; \theta)) d\mu(z)$ à maximiser. C'est souvent ce qui se passe en pratique.

• Pourquoi EM plutôt que GS? Pour éviter des calculs d'espérance du gradient; on va voir comment tout de suite.

Terminologie: y_1, \dots, y_n : données observées tq $\forall i \exists z_i$ tq $y_i | z_i \sim q(y_i | z_i; \theta)$
 z_1, \dots, z_n : données manquantes
 (\vec{y}, \vec{z}) : données complètes.

Rq: • écrit sous la forme (*), c'est un problème de maximum de vraisemblance.

Si θ est considéré comme une v.a. sur laquelle on a un a priori, on peut aussi utiliser ce cadre de l'EM pour un maximum a posteriori

• Si on doit maximiser une intégrale DÉTERMINISTE, en réécrivant \int en $\mathbb{E}_{\mu} \rightarrow \mathbb{E}_{\theta}$.

1) D'où vient l'EM?

- * Pour des raisons de commodité, on se place dans le cadre éventuellement Bayésien (\rightarrow)
- * On commence avec un θ initial. L'état courant de la suite est noté θ^t .
- * Trouver $\underset{\theta \in \Theta}{\text{argmax}} q(y, \theta)$ est équivalent à trouver $\theta^* = \text{argmax} \log q(y, \theta)$
Le problème est qu'il est difficile de trouver cet argmax du fait de l'J.
- * idée: Trouver $B(\theta, \theta^t) \leq \log q(\theta | y)$ et de maximiser cette fonction B au lieu de $q(\cdot | y)$ \rightarrow en itérant, on voit que l'on peut espérer une convergence vers un max local puisque B est "optimale" à chaque étape.
- * On réécrit le problème pour trouver cette fonction B optimale.

$$\log q(y, \theta) = \log \int_{\mathbb{Z}} q(y, z, \theta) d(z) = \log \int_{\mathbb{Z}} f^t(z) \frac{q(y, z, \theta)}{f^t(z)} d(z)$$

où $f^t(z)$ est une densité non nulle sur \mathbb{Z} dépendant de la valeur courante de θ .

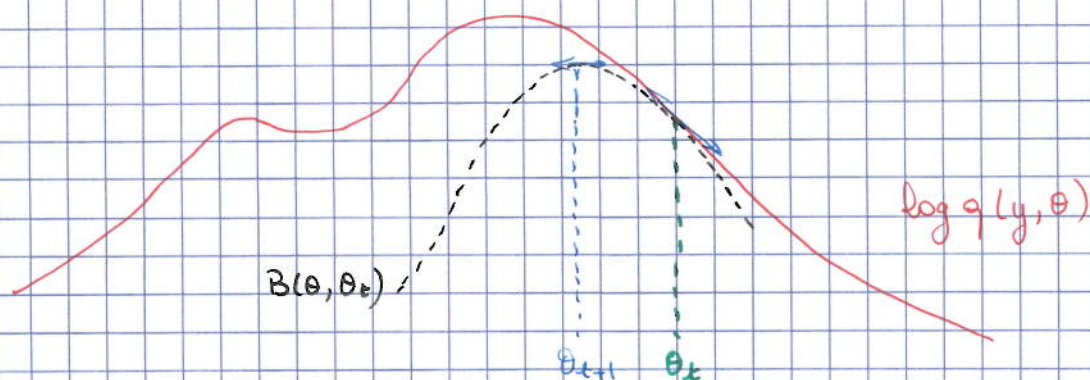
Par l'inégalité de Jensen: $\log \int_{\mathbb{Z}} f^t(z) \frac{q(y, z, \theta)}{f^t(z)} dz \geq \int_{\mathbb{Z}} f^t(z) \log \left[\frac{q(y, z, \theta)}{f^t(z)} \right] dz$

Cette inégalité de convexité est optimale (car \exists cas d'égalité) donc on

pose $B(\theta, \theta^t) = \int_{\mathbb{Z}} f^t(z) \log \left[\frac{q(y, z, \theta)}{f^t(z)} \right] dz$

- * Maintenant que l'on a la forme générale de B , il faut trouver la meilleure densité $f^t(z)$ qui permet que l'inégalité soit la meilleure possible.

L'idée est de prendre $f^t(z)$ tq B touche la fonction objectif $\log q(y, \theta)$ en l'état courant θ^t :



En équation : $B(\theta, \theta_t) = \int f_t(z) \log \left[\frac{q(y, z, \theta)}{f_t(z)} \right] dz$ (3)

On veut que $B(\theta_t, \theta_t)$ touche $q(y, \theta_t)$ donc comme B est un minorant on optimise en f_t en maximisant $B(\theta_t, \theta_t)$ par dérivation en f_t .

On cherche donc $f_t^* = \arg \max_{f_t} \int f_t(z) \log \left[\frac{q(y, z, \theta)}{f_t(z)} \right] dz$ sous contrainte

de densité : $f_t(z) \geq 0 \forall z$ et $\int f_t(z) dz = 1$

On introduit un multiplicateur de Lagrange :

$$\begin{aligned} G(f_t) &= \lambda \left[1 - \int f_t(z) dz \right] + \int f_t(z) \log \left[\frac{q(y, z, \theta)}{f_t(z)} \right] dz \\ &= \lambda \left[1 - \int f_t(z) dz \right] + \int f_t(z) \log q(y, z, \theta) dz - \int f_t(z) \log f_t(z) dz \end{aligned}$$

$$\begin{cases} \frac{\partial G}{\partial f_t(z)} = -\lambda + \log q(y, z, \theta) - \log f_t(z) - 1 = 0 & (1) \end{cases}$$

$$\begin{cases} \frac{\partial G}{\partial \lambda} = 1 - \int f_t(z) dz & (2) \end{cases}$$

$$(1) \Leftrightarrow \log f_t(z) = \log q(y, z, \theta) - \lambda - 1$$

$$(2) \Leftrightarrow 1 = \int \exp \left[(-\lambda - 1) q(y, z, \theta) \right] dz = e^{-\lambda - 1} \int q(y, z, \theta) dz$$

$\stackrel{\text{marginale}}{=} e^{-\lambda - 1} \int q(y, \theta)$

$$\text{D'où } \boxed{f_t(z) = q(y, z, \theta) \times \frac{1}{q(y, \theta)} = q(z|y, \theta)}$$

On vérifie facilement que $B(\theta_t, \theta_t) = q(y, \theta_t)$:

$$\begin{aligned} B(\theta_t, \theta_t) &= \int q(z|y, \theta_t) \log \left[\frac{q(y, z, \theta_t)}{q(z|y, \theta_t)} \right] dz \\ &= \int q(y, z, \theta_t) \frac{1}{q(y, \theta_t)} \log \left[\frac{q(y, z, \theta_t) q(y, \theta_t)}{q(y, z, \theta_t)} \right] dz \\ &= \frac{\log q(y, \theta_t)}{q(y, \theta_t)} \int q(y, z, \theta_t) dz = \log q(y, \theta_t) \times \frac{q(y, \theta_t)}{q(y, \theta_t)} \\ &= \log q(y, \theta_t) \end{aligned}$$

* Dernière étape : il faut maximiser $B(\theta, \theta_t)$ en θ .

$$B(\theta, \theta_t) = \mathbb{E}_{f_t(z)} [\log q(y, z, \theta) - \log f_t(z)] \quad (4)$$

$$= \mathbb{E}_{f_t} [\log q(y, z, \theta) + \log q(\theta)] + \mathcal{H}(f_t)$$

$\xrightarrow{\text{la loi de la cas Bayésien}} \quad \xrightarrow{\text{entropie de } f}$

$$= Q(\theta | \theta_t) + \log q(\theta) + \underbrace{\mathcal{H}(f_t)}_{\text{indépendante de } \theta \text{ à optimiser}}$$

Donc $\theta_{t+1} = \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta | \theta_t) (+ \log q(\theta) \text{ la loi cas Bayésien})$

2) L'algorithme en deux étapes :

(E): Calculer $Q(\theta | \theta_t) = \mathbb{E}[\log q(y, z, \theta) | y^*, \theta_t]$

où l'espérance est prise contre la loi a posteriori des données manquantes en le paramètre courant.

(M) $\theta_{t+1} = \underset{\theta \in \Theta}{\operatorname{argmax}} Q(\theta | \theta_t) (+ \log q(\theta))$

Remarque: Ne nécessite plus de calculer qu'une E et non J et DJS des GS.
Cet avantage reste un inconvénient : il faut calculer cette J sous la loi a posteriori...

3) Exemple:

L'exemple le plus courant est celui d'un modèle de mélange :

la densité des observations est $f_\theta(y) = \sum_{j=1}^m \alpha_j \cdot f_j(y)$

le modèle sous-jacent est celui-ci : $\left\{ \begin{array}{l} P(Z_k = j) = \alpha_j : Z_k \sim \sum_{j=1}^m \alpha_j \delta_j \\ (**) \quad Y_k | Z_k = j \sim f_j(y) \end{array} \right.$

On cherche $\hat{\theta}_n = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{j=1}^m \log f_\theta(Y_j)$ où θ sont les paramètres de f_j et aussi les poids α_i

Supposons ici que $f_\theta(y) = \sum_{j=1}^m \alpha_j \frac{\lambda_j^y}{y!} e^{-\lambda_j}$: mélange de lois de Poisson

la vraisemblance complète (en log) est donc :

$$\log q_\theta(Y, Z) = -\log(Y!) + \sum_{j=1}^m [\log \alpha_j - \lambda_j] \mathbb{1}_{\{Z=j\}} + \sum_{j=1}^m \log(\lambda_j) Y \mathbb{1}_{\{Z=j\}}$$

$$Q(\theta|\theta_t) = E \left[\log q_\theta(Y, Z) \mid Y, \theta_t \right]$$

où $q(Z|Y, \theta_t)$ est une loi discrète.

$$Q(\theta|\theta_t) = \int \sum_{i=1}^n \left[\log(Y_i!) + \sum_{j=1}^m [\log \alpha_j - \lambda_j] \mathbb{1}_{\{Z_i=j\}} + \sum_{j=1}^m \log(\lambda_j) Y_i \mathbb{1}_{\{Z_i=j\}} \right] q(Z_i | Y_i, \theta_t) dZ_i$$

$$= - \sum_{i=1}^n \log(Y_i!) + \sum_{i=1}^n \sum_{j=1}^m (\log \alpha_j - \lambda_j) P(Z_i=j | Y_i, \theta_t)$$

$$+ \sum_{i=1}^n \sum_{j=1}^m \log(\lambda_j) Y_i P(Z_i=j | Y_i, \theta_t)$$

$$= C_{\text{indépendante de } Z, \theta} + \left[\sum_{j=1}^m (\log \alpha_j - \lambda_j) \times \left[\sum_{i=1}^n P(Z_i=j | Y_i, \theta_t) \right] \right] + \left[\sum_{j=1}^m \log(\lambda_j) \times \sum_{i=1}^n Y_i P(Z_i=j | Y_i, \theta_t) \right]$$

Maintenant la mise à jour de λ_j^{t+1} et α_j^{t+1} via : par dérivation :

- α_j^{t+1} est solution d'une minimisation sous contrainte; comme précédemment avec un multiplicateur de Lagrange $\mu \times (\sum_{j=1}^m \alpha_j - 1)$ on trouve :

$$\alpha_j^{t+1} = \frac{1}{n} \sum_{i=1}^n P(Z_i=j | Y_i, \theta_t)$$

$$\lambda_j^{t+1} = \frac{\sum_{i=1}^n Y_i P(Z_i=j | Y_i, \theta_t)}{\sum_{i=1}^n P(Z_i=j | Y_i, \theta_t)}$$

Le calcul de $P(Z_i=j | Y_i, \theta_t)$ se fait à partir du modèle généralisé (**).

4) Convergence:

Une remarque tout d'abord: La simplicité apparente de cet algorithme dans le cas de notre exemple n'est pas le fait de la simplicité de ce modèle en particulier. Elle vient du fait que l'on peut écrire simplement la vraisemblance complète $q(Y, Z; \theta_t)$ et que son log donne une fonction simple qui sépare des termes dépendant de Z et ceux de θ .

C'est une propriété que l'on retrouve comme condition de convergence.

Convergence montrée (puis corrigée) par Dempster, Laird, Rubin et Wu (1983)
puis reformulée avec des hypothèses plus simples par Delyon, Lavielle,
Moulines (1999)

(M1) L'espace des paramètres Θ est un ouvert de \mathbb{R}^p et la vraisemblance complète du modèle est donnée par : $q(y, z; \theta) = \exp[-\psi(\theta) + \langle S(y, z), \psi(\theta) \rangle]$, où $\langle ; \rangle$ est le produit scalaire euclidien dans \mathbb{R}^m , S est une fonction borélienne en z sur \mathbb{R}^p à valeurs dans $\mathcal{Y} \subset \mathbb{R}^m$. S est appelée statistique exhaustive du modèle (unique à C^0 près)

On suppose de plus que l'enveloppe convexe de $S(\mathbb{R}^p) \subset \mathcal{Y}$ et que $\forall \theta \in \Theta$

$$\int_{\mathbb{R}^p} |S(y, z)| q(z|y, \theta) dz < \infty$$

(M2) ψ et ψ sont C^1 sur Θ

(M3) $\bar{S} : \Theta \rightarrow \mathcal{Y}$ tq $\bar{S}(\theta) \triangleq \int_{\mathbb{R}^p} S(y, z) q(z|y, \theta) dz$ est C^1 sur Θ

(M4) $\ell(\theta) \triangleq \log q(y, \theta)$ est C^1 sur Θ et

$$\partial_\theta \int q(y, z, \theta) dz = \int \partial_\theta q(y, z, \theta) dz$$

(M5) Soit $L : \mathcal{Y} \times \Theta \rightarrow \mathbb{R}$ tq $L(x, \theta) \triangleq -\psi(\theta) + \langle S(y, z), \psi(\theta) \rangle$

Alors $\exists \hat{\theta} : \mathcal{Y} \rightarrow \Theta$ C^1 sur \mathcal{Y} tq

$$\forall x \in \mathcal{Y} \quad \forall \theta \in \Theta \quad L(x, \hat{\theta}(x)) \geq L(x, \theta)$$

on définit $\mathcal{A} = \{\theta \in \Theta \text{ tq } \partial_\theta \ell(\theta) = 0\}$ et $d(x, \mathcal{A})$ la distance euclidienne du point x à l'ensemble fermé \mathcal{A} .

Théorème (D, L, M, 99) On suppose que le modèle vérifie (M1-5).

Alors $\forall \theta_0 \in \Theta$, la suite $(\ell(\theta_k))_k$ donnée par l'algorithme EM est croissante et la suite $(\theta_k)_k$ vérifie $\lim_{k \rightarrow \infty} d(\theta_k, \mathcal{A}) = 0$

On ne va pas prouver ici ce théorème. Juste commenter les hypothèses. Je renvoie au D47 pour les détails.

Rq: • $(\theta_k)_k$ est une suite déterministe donc elle dépend de θ_0 . Il se peut que pour un θ_0 on obtienne un max local \Rightarrow tester pour plusieurs θ_0 et choisir celui pour lequel $l(\theta_k)$ maximal.

• L'écriture sous forme exponentielle: beaucoup de modèles entre dans cette catégorie même des très complexes. \rightarrow pas tellement une contrainte.

Dans notre exemple:

$$S_{1,j}(y_i, z_i) = \mathbb{1}_{\{z_i=j\}}$$

$$S_{2,j}(y_i, z_i) = y_i \mathbb{1}_{\{z_i=j\}}$$

\rightarrow S de dimension 2m

$$\psi_j(\theta) = \log(\alpha_j) - \lambda_j$$

$$\psi_{2,j}(\theta) = \log \lambda_j$$

\rightarrow en exo

• (H4) garantit une certaine régularité de log vraisemblance

• (H5) permet une mise à jour facile de θ connaissant S.

Dans notre exemple:

$$\hat{\alpha}_j^{t+1} = \frac{1}{n} \sum_{i=1}^n S_{1,i}(y_i, z_i)$$

$$\hat{\lambda}_j^{t+1} = \frac{\sum_{i=1}^n S_{2,i}(y_i, z_i)}{\sum_{i=1}^n S_{1,i}(y_i, z_i)}$$

en exo

III Variantes de l'EM:

1) SEM: (Simulated EM) **Geyer, Diebolt (1986)**

• cherche à éviter le calcul de l'espérance selon la loi a posteriori

• On simule $z_i^k \sim q(\cdot | y_i, \theta_k)$

• On maximise $\sum_{i=1}^n \log q(y_i, z_i^k, \theta)$

⊕ \rightarrow évite le calcul de E

\rightarrow l'aliat introduit réduit la dépendance en θ_0 car exploration plus générale des modes de $q(z|y, \theta)$

⊖ \rightarrow Convergence prouvée en moyenne seulement (ps mélange à du recuit)

2) MCEM: (Monte Carlo EM) Wei, Tanner (1990)

- Consiste à approcher l'espérance conditionnelle par une somme de Monte Carlo.
- On simule un (grand) nombre T de variables $(z_i^t)_{t \in \{1, \dots, T\}} \sim q(z|y_i, \theta_k)$ iid
- On approche $Q(\theta|\theta_k) \approx \frac{1}{T} \sum_{t=1}^T \log q(y_i, z_i^t, \theta_k)$

⊕ Allongement du temps de calcul

⊕ Pas de résultats théoriques de convergence connus.

3) SAEM: Stochastic Approximation - EM: Delyon (1999)

L'idée est de profiter du passé des simulations avec $\theta_0, \dots, \theta_{k-1}$ pour mettre à jour notre approximation de $Q(\theta|\theta_k)$. Cela est fait à l'aide des approximations stochastiques.

L'étape E est divisée ici en deux sous étapes :

S: simulation de $z_k \sim q(z|y, \theta_k)$ (pour tout i , on l'omet pour des raisons de commodité)

A: Approximation stochastique: soit $(\gamma_k)_k$ une suite de pas:

$$\begin{cases} Q_{k+1}(\theta) = Q_k(\theta) + \gamma_k (\log q(y, z_k, \theta) - Q_k(\theta)) \\ \theta_0 \text{ donné} \end{cases}$$

π : maximisation: $\theta_{k+1} = \underset{\theta \in \Theta}{\operatorname{argmax}} Q_{k+1}(\theta)$

⊕ Ne demande qu'une simulation à chaque étape

• garde en mémoire le passé via la suite $(\gamma_k)_k$

• Si le modèle est exponentiel et $\exists \hat{\theta}$ comme dans π_1 et π_5 alors l'AS se fait uniquement sur S ce qui produit une suite $(\theta_k)_k$ tq

$$\theta_{k+1} = \theta_k + \gamma_k (S(y, z_k) - \theta_k)$$

$$\text{et } \theta_{k+1} = \hat{\theta}(\theta_{k+1})$$

$$\underline{\text{Rq:}} \quad Q_{k+1} = (1 - \gamma_k) Q_k + \gamma_k \log q(y, z_k, \theta_k)$$

$$\begin{aligned} &= (1 - \gamma_k) \left((1 - \gamma_{k-1}) Q_{k-1} + \gamma_{k-1} \log q(y, z_{k-1}, \theta_{k-1}) \right) + \gamma_k \log q(y, z_k, \theta_k) \\ &= \sum_{p=1}^k \mu_p \log q(y, z_p, \theta_p) \quad \text{avec } \sum_{p=1}^k \mu_p = 1 \Rightarrow \text{Approx de type MC} \end{aligned}$$

• Résultats de convergence:

(9)

En plus de (H1-5) il faut des conditions supplémentaires qu'on dérive en regardant la somme précédente. Elles portent sur $(Y_k)_k$ et la continuité de l et $\hat{\theta}$.

Sont $(\mathcal{F}_k)_k$ la filtration engendrée par la suite $(z_k)_k$ des simulations.

(SAEN1) $\forall k \in \mathbb{N} \quad Y_k \in [0,1], \quad \sum_{k \geq 0} Y_k = +\infty, \quad \sum_{k \geq 0} Y_k^2 < \infty$

(SAEN2) $l: \Theta \rightarrow \mathbb{R}$ et $\hat{\theta}: \mathcal{Y} \rightarrow \Theta$ sont m fois continuellement différentiables où m est tel que \mathcal{Y} ouvert de \mathbb{R}^m .

(SAEN3) • $\forall \phi$ bornée positive

$$E[\phi(z_{k+1}) | \mathcal{F}_k] = \int \phi(z) q(z | y, \theta_k) dz$$

$$\bullet \forall \theta \in \Theta \quad \int \|S(y, z)\|^2 q(z | y, \theta) dz < \infty \text{ et}$$

$$r(\theta) \triangleq \text{Var}(S(y, z) | y, \theta) \triangleq \int_{\mathbb{R}^d} S(y, z)^2 q(z | y, \theta) dz - \left[\int_{\mathbb{R}^d} S(y, z) q(z | y, \theta) dz \right]^2$$

est C^0 en Θ

Théorème (DLN) On suppose (H1-5) et (SAEN1-3) vérifiées ainsi que

(C) $(\theta_k)_{k \geq 0}$ est à valeur dans un ensemble compact de \mathcal{Y}

alors $\lim_{k \rightarrow \infty} d(\theta_k, \mathcal{L}) = 0$ ps.

Rq: (C) intègre de convergence

Θ vers un point stationnaire de $l \dots$ mais pas point selle;

il faut des conditions type convexité au moins locale pour conclure à un maximum.

Sous des conditions un peu plus restrictives qui ressemblent à celles vues pour les A.S, on obtient aussi un TCL.

→ détails DLN 99