

**Consignes :** Vous devez séparer pour chacun des exercices votre programme principal (l'équivalent du main) du reste du code (fonctions etc). Le dépôt de votre travail se fera, sous la forme d'une archive, à déposer sur ecampus (exclusivement), le 20/11/2018 avant 12h30, heure stricte. (Les envois par mail ne seront pas pris en compte).

Le barème (4pts Etape E, 8 pts étape M, 3pts le reste du code) est donné à titre indicatif. Les exercices pour lesquels les codes ne s'exécutent pas, seront notés sur la moitié des points du barème.

On considère une série temporelle  $(Y_1, \dots, Y_n)$ , avec  $Y_t \in \mathbb{N}$ , régie par un processus latent  $(Z_1, \dots, Z_n)$  à  $K$  états où  $Z_t$  représente l'état latent à l'instant  $t$  ( $t = 1, \dots, n$ ).  $Z_t \in \{1, 2, \dots, K\}$ ,  $K \in \mathbb{N}^*$ .

On suppose que pour chaque état  $k$ , la variable observée  $Y_t$  est modélisé en fonction du temps par le modèle log-linéaire suivant

$$\ln(\lambda_k(t, \boldsymbol{\beta}_k)) = \boldsymbol{\beta}_k^T \mathbf{x}_t \quad (1)$$

où  $\mathbf{x}_t = (1, t)^\top$  et  $Y_t | Z_t = k; \boldsymbol{\beta}_k \sim \mathcal{P}(\lambda_k(t; \boldsymbol{\beta}_k))$  de loi Poisson

$$\mathbb{P}(Y_t = y_t | Z_t = k; \lambda_k(t; \boldsymbol{\beta}_k)) = \frac{e^{-\lambda_k(t; \boldsymbol{\beta}_k)} (\lambda_k(t; \boldsymbol{\beta}_k))^{y_t}}{y_t!}, \forall y_t \in \mathbb{N}, \quad (2)$$

qu'on notera  $\text{Poisson}(y_t; \lambda_k(t, \boldsymbol{\beta}_k))$ . La probabilité de l'état  $k$  à l'instant  $t$  définie par le modèle logistique :

$$\mathbb{P}(Z_t = k; \mathbf{w}) = \frac{e^{\mathbf{w}_k^\top \mathbf{x}_t}}{1 + \sum_{\ell=1}^K e^{\mathbf{w}_\ell^\top \mathbf{x}_t}}, \quad (3)$$

où  $\mathbf{w} = (\mathbf{w}_1^\top, \dots, \mathbf{w}_{K-1}^\top)^\top$  et  $\mathbf{w}_K = \mathbf{0}$ . On notera par  $\pi_k(t; \mathbf{w})$  cette probabilité.

La loi de la variable observée est dans ce cas définie par :

$$\mathbb{P}(Y_t = y_t; \boldsymbol{\theta}) = \sum_{k=1}^K \pi_k(t; \mathbf{w}) \text{Poisson}(y_t; \lambda_k(t, \boldsymbol{\beta}_k)) \quad (4)$$

$\boldsymbol{\theta}$  étant le vecteur paramètre inconnu du modèle.

L'objectif est de segmenter la série temporelle sur la base du modèle (4). Pour cela, on estime  $\boldsymbol{\theta}$  à partir des données en maximisant la log-vraisemblance  $L(\boldsymbol{\theta}) = \sum_{t=1}^n \ln \mathbb{P}(Y_t = y_t; \boldsymbol{\theta})$  par l'algorithme EM. On montre que celui-ci consiste à partir d'un modèle initial de paramètre  $\boldsymbol{\theta}^{(0)}$  et alterner à chaque itération  $q$  entre les deux étapes E- et M- suivantes jusqu'à la convergence :

**Étape E :** Calculer la fonction  $Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)})$  définie par :

$$Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)}) = \sum_{t=1}^n \sum_{k=1}^K \tau_k(t; \boldsymbol{\theta}^{(q)}) \ln[\pi_k(t; \mathbf{w}) \text{Poisson}(y_t; \lambda_k(t, \boldsymbol{\beta}_k))] \quad (5)$$

où  $\tau_k(t; \boldsymbol{\theta}^{(q)}) = \mathbb{P}(Z_t = k | y_t; \boldsymbol{\theta}^{(q)})$  est la probabilité a posteriori de l'état  $k$  à l'instant  $t$ , définie par :

$$\tau_k(t; \boldsymbol{\theta}^{(q)}) = \frac{\pi_k(t; \mathbf{w}^{(q)}) \text{Poisson}(y_t; \lambda_k(t, \boldsymbol{\beta}_k^{(q)}))}{\mathbb{P}(Y_t = y_t; \boldsymbol{\theta}^{(q)})}, \quad \forall t = 1, \dots, n \quad (6)$$

**Étape M :** Mettre à jour les paramètres du modèle en calculant  $\boldsymbol{\theta}^{(q+1)}$  définie par :

$$\boldsymbol{\theta}^{(q+1)} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}; \boldsymbol{\theta}^{(q)}). \quad (7)$$

Ce problème de maximization n'admet pas de solution analytique.

1. La maximisation de (5) par rapport aux vecteurs  $\boldsymbol{\beta}_k$ , peut se faire par l'algorithme Newton-Raphson.

A chaque itération  $m$  de cet algorithme, la mise à jour est donnée par l'équation suivante :

$$\boldsymbol{\beta}_k^{(m+1)} = \boldsymbol{\beta}_k^{(m)} + \left[ \sum_{t=1}^n \tau_k(t; \boldsymbol{\theta}^{(q)}) e^{\boldsymbol{\beta}_k^{\top} \mathbf{x}_t} \mathbf{x}_t \mathbf{x}_t^{\top} \right]_{\boldsymbol{\beta}_k = \boldsymbol{\beta}_k^{(m)}}^{-1} \sum_{t=1}^n \tau_k(t; \boldsymbol{\theta}^{(q)}) \mathbf{x}_t \left( y_t - e^{\boldsymbol{\beta}_k^{\top} \mathbf{x}_t} \right)_{\boldsymbol{\beta}_k = \boldsymbol{\beta}_k^{(m)}} \quad (8)$$

2. On prend  $K = 2$ . Dans ce cas on montre que la mise à jour itérative du vecteur  $\mathbf{w}$  à chaque itération  $s$  de l'algorithme Newton-Raphson est donnée par :

$$\mathbf{w}^{(s+1)} = \mathbf{w}^{(s)} + \left[ \sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^{\top} \pi(t; \mathbf{w}) (1 - \pi(t; \mathbf{w})) \right]_{\mathbf{w} = \mathbf{w}^{(s)}}^{-1} \sum_{t=1}^n \mathbf{x}_t \left( \tau_t(\boldsymbol{\theta}^{(q)}) - \pi(t; \mathbf{w}) \right)_{\mathbf{w} = \mathbf{w}^{(s)}} \quad (9)$$

où  $\pi(t; \mathbf{w}) = \frac{e^{\mathbf{w}^{\top} \mathbf{x}_t}}{1 + e^{\mathbf{w}^{\top} \mathbf{x}_t}}$  est la probabilité de l'un des deux états.

3. Soit  $\hat{\boldsymbol{\theta}}$  le vecteur paramètre estimé par l'algorithme EM. On peut en déduire la séquence des états  $(\hat{Z}_1, \dots, \hat{Z}_n)$  où

$$\hat{Z}_t = \arg \max_k \pi_k(t; \hat{\mathbf{w}}). \quad (10)$$

**Travail demandé :** Implémenter l'algorithme EM et appliquer le à la série dataPoisson.mat avec  $K = 2$ .

Il est conseillé d'utiliser une formulation vectorielle des mises à jour des paramètres  $\boldsymbol{\beta}_k$  et  $\mathbf{w}$ .

- Afficher dans un graphique la série temporelle
- Afficher dans un graphique les probabilités logistiques estimées
- Afficher dans un graphique la log-vraisemblance calculée à chaque itération de l'algorithme EM
- Afficher dans un graphique la série temporelle, et la fonction de régression de Poisson de chaque état  $k$  définie par :  $\lambda_k(t, \hat{\boldsymbol{\beta}}_k) = \mathbb{E}[Y_t | Z_t = k; \hat{\boldsymbol{\beta}}_k] = e^{\hat{\boldsymbol{\beta}}_k^{\top} \mathbf{x}_t} \quad \forall t = 1, \dots, n$ .