

# Estimation des paramètres d'une distribution Gamma

**Préparation** : avant la séance répondez aux questions 4, 5a, 5c, 5d

L'ESTIMATION de paramètres est une question cruciale pour le traitement de données expérimentales. A titre d'exemple, on s'intéresse ici à un « cas d'école », relativement simple concernant l'estimation des deux paramètres d'une distribution Gamma. L'objet de ce travail n'est pas une étude de cette distribution (moyenne, variance, ... densité, fonction de répartition, ... propriétés diverses, ...) mais l'identification (l'estimation) de ses paramètres.

**Fonction gamma** – On rappelle la définition de la fonction  $\Gamma$  :

$$\Gamma(u) = \int_{\mathbb{R}_+} t^{u-1} \exp -t \, dt, \text{ pour } u \in \mathbb{R}_+.$$

Cette fonction interviendra essentiellement au travers de son logarithme  $\log \Gamma$  et la routine *Matlab* `gamma1n` permet son calcul. Sa dérivée est la fonction di-Gamma notée  $\psi$

$$\psi(u) = \frac{d}{du} \log \Gamma(u),$$

disponible par la routine *Matlab* `psi`. Cette même routine permet de calculer les dérivées de  $\psi$ . On fournit en plus la routine `InvPsi` qui permet d'inverser  $\psi$ .  $\triangle$

On s'intéresse à la densité Gamma paramétrée par  $\alpha > 0$  et  $\beta > 0$  :

$$f_{X|A,B}(x|\alpha, \beta) = \frac{\beta^\alpha}{\Gamma[\alpha]} x^{\alpha-1} \exp [-\beta x] \mathbb{1}_+(x), \quad (1)$$

où  $\mathbb{1}_+$  est l'indicatrice de  $\mathbb{R}_+$ . On rappelle que sa moyenne et sa variance valent respectivement :  $\alpha/\beta$  et  $\alpha/\beta^2$ . Lorsque  $\alpha > 1$ , son maximiseur est  $(\alpha - 1)/\beta$ .

On en observe  $N$  échantillons  $x_1, x_2, \dots, x_N$  indépendants et on les regroupe dans le vecteur  $\mathbf{x} \in \mathbb{R}_+^N$ . Pour alléger certaines écritures, on note

$$\begin{aligned} M_A(\mathbf{x}) &= \frac{1}{N} \sum x_n & \text{et} & & V_A(\mathbf{x}) &= \frac{1}{N} \sum (x_n - M_A(\mathbf{x}))^2 \\ M_L(\mathbf{x}) &= \frac{1}{N} \sum \log x_n & \text{et} & & V_L(\mathbf{x}) &= \frac{1}{N} \sum (\log x_n - M_L(\mathbf{x}))^2 \end{aligned}$$

les moyennes et variances empiriques des données et de leur logarithme.

1. « **Théorie** » — En utilisant la routine *GammaPDF*, tracez la densité de probabilité pour diverses valeurs de  $(\alpha, \beta)$ . Attention, *Matlab* fonctionne à l'inverse pour le second paramètre : il faut lui passer  $(\alpha, 1/\beta)$  pour lui causer de la densité définie par (1).

**2. « Simulations »** — Tirez les  $N$  réalisations de la variable en utilisant la routine *GammaRND*. Tracez-les sous la forme d'un histogramme que vous positionnerez en regard la densité.

**3. Lien « empirique – vrai ».**

**3a.** Sous Matlab, calculez les moments (moyenne et variance) empiriques et vérifiez qu'ils prennent la « bonne » valeur compte-tenu des vraies valeurs des paramètres  $\alpha$  et  $\beta$ .

**3b.** A l'inverse, en approchant les deux moments vrais (fonction des vrais paramètres) par les deux moments empiriques (calculés à partir des données), proposez un estimateur empirique, simple et rapide, pour  $\alpha$  et  $\beta$ . On parle d'estimation par une méthode de moments.

**4.** Écrivez la vraisemblance du couple  $(\alpha, \beta)$  attachée aux données  $\mathbf{x}$  ainsi que la log-vraisemblance que l'on notera  $LV_{\mathbf{x}}(\alpha, \beta)$ .

On se donne une densité *a priori*  $\pi_{A,B}(\alpha, \beta)$  pour le couple  $(A, B)$  que l'on choisit uniforme sur le pavé  $[0, \alpha_M] \times [0, \beta_M]$ , où  $\alpha_M$  et  $\beta_M$  sont des valeurs maximales prédéfinies pour les paramètres.

$$\pi_{A,B}(\alpha, \beta) = \begin{cases} (\alpha_M \beta_M)^{-1} & \text{si } 0 < \alpha < \alpha_M \text{ et } 0 < \beta < \beta_M \\ 0 & \text{sinon} \end{cases}$$

On pourra prendre par exemple  $\alpha_M = 8$  et  $\beta_M = 5$ . Concernant la vraie valeur des paramètres, commencera, par exemple, avec  $\alpha^* = 3$  et  $\beta^* = 2$ .

**5. Travail sur la distribution *A posteriori***

**5a.** Déterminez la densité *a posteriori*  $\pi_{A,B|\mathbf{X}}(\alpha, \beta|\mathbf{x})$  à un facteur multiplicatif près et surtout son logarithme, noté  $LVP_{\mathbf{x}}(\alpha, \beta)$ , à un terme additif près.

**5b.** Tracez  $LVP_{\mathbf{x}}(\alpha, \beta)$  sur une grille de valeur du pavé  $[0, \alpha_M] \times [0, \beta_M]$ . La fonction *meshgrid* pourra être utile.

**5c.** Déduisez-en le logarithme de chacune des conditionnelles *a posteriori*  $\pi_{A|\mathbf{X},B}(\alpha|\mathbf{x}, \beta)$  et  $\pi_{B|\mathbf{X},A}(\beta|\mathbf{x}, \alpha)$  que l'on notera  $LVP_{\mathbf{x}}(\alpha|\beta)$  et  $LVP_{\mathbf{x}}(\beta|\alpha)$ .

**5d.** Calculez le maximiseur de chaque conditionnelle *a posteriori*. Naturellement, le maximiseur par rapport à  $\alpha$  sera fonction de  $\beta$  et inversement. On pourra faire intervenir l'inverse de la fonction  $\psi$  évoquée en introduction.

On déduit de ce qui précède, un algorithme de calcul du maximiseur *a posteriori*. On montre qu'on peut maximiser la densité *a posteriori* par rapport à  $(\alpha, \beta)$  en maximisant alternativement vis-à-vis de chacune des deux variables  $\alpha$  et  $\beta$ , c'est-à-dire en maximisant à tour de rôle chaque conditionnelle  $(\alpha|\mathbf{x}, \beta)$  et  $(\beta|\mathbf{x}, \alpha)$ . Une forme générique de l'algorithme d'optimisation est alors la suivante.

1. Initialisation :  $\beta^{[1]} = 1$
2. Boucle pour  $k = 2, \dots, K$ 
  - Déterminez  $\alpha^{[k]}$  comme maximiseur de  $\text{LVP}_x(\alpha|\beta^{[k-1]})$
  - Déterminez  $\beta^{[k]}$  comme maximiseur de  $\text{LVP}_x(\beta|\alpha^{[k]})$
3. Fin

Dans le langage *Matlab*, le codage prend la forme suivante dont vous pouvez vous inspirer et qu'il vous appartient de compléter.

```
% Initialisation
BetaCourant = 1;

% Boucle d'optimisation variable par variable
while Delta ...

    % Mises A Jour Alpha
    AlphaCourant = ...
    GardeAlpha = [GardeAlpha AlphaCourant];
    GardeBeta = [GardeBeta BetaCourant];

    % Mises A Jour Beta
    BetaCourant = ...
    GardeAlpha = [GardeAlpha AlphaCourant];
    GardeBeta = [GardeBeta BetaCourant];

    % Variations
    Delta = ...

end
% Fin de boucle d'optimisation
```

- 6a.** Représentez les deux ensembles d'itérées produites  $\alpha^{[k]}$  et  $\beta^{[k]}$  pour  $k = 1, \dots, K$  en fonction de  $k$  sur deux sous figures d'une même figure.
- 6b.** Représentez également, dans le plan, les couples d'itérées  $(\alpha^{[k]}, \beta^{[k]})$  superposés à la représentation de la densité *a posteriori* de la question 5a.
- 7.** Récupérez le fichier fourni *Mystere.mat* : il contient un jeu de données résultant de valeurs de paramètres  $\alpha$  et  $\beta$  inconnues de vous. Déterminez les valeurs de ces paramètres.