

Vers la découverte et la sélection : une approche bayésienne

— Projet BHI-PRO —

Cette proposition concerne la question de l'identification de bio-marqueurs protéiques dans le cadre du projet BHI-PRO soutenu par l'ANR. En résumé, on s'intéresse à une situation dans laquelle on observe les concentrations pour un jeu de protéines relatif à deux populations d'individu et on se pose la question de la possibilité de discriminer entre les deux populations sur la base de ces protéines.

Pour cela, on considère une population d'individus indexés par n et on s'intéresse à leur statut biologique \mathbf{b}_n qui prend deux valeurs \mathbf{m} (malade) ou \mathbf{s} (sain) et à leurs concentrations en P protéines regroupées dans $\mathbf{x}_n \in \mathbb{R}^P$. On observe N individus et on note $\mathbf{b} = [\mathbf{b}_1, \dots, \mathbf{b}_N]$ le vecteur des N statuts et $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ la matrice des N vecteurs de P concentrations. On note $\mathcal{I}_{\mathbf{b}}$ et $N_{\mathbf{b}}$ l'ensemble des indices des individus de statut \mathbf{b} et son cardinal.

Le statut est décrit par une variable de Bernoulli (\mathcal{B}) de paramètre p et les concentrations sont décrites (conditionnellement au statut) par des densités gaussiennes (\mathcal{N}) de moyennes et précisions $(\mathbf{m}_{\times}, \Gamma_{\times})$. Ces paramètres sont également probabilisés : pour des questions pratiques, on choisit un modèle Bêta (\mathcal{P}) pour p et des modèles Gauss-Wishart (\mathcal{NW}) pour les $(\mathbf{m}_{\times}, \Gamma_{\times})$. Les hyperparamètres sont fixés de manière à inclure les informations disponibles sur les paramètres. Dans le présent travail, on se focalise sur les a priori diffus (non-informatif).

On se pose la question¹ du caractère discriminant du jeu de protéines considéré : « décider » si ce jeu témoigne de l'existence de deux sous-populations distinctes ou bien en réalité d'une seule. Dans le premier cas on dira que le modèle est discriminant (on écrira $\mathfrak{M} = \oplus$) et dans le second on dira qu'il ne l'est pas (on écrira $\mathfrak{M} = \ominus$).

Deux modèles concurrents

- $\mathfrak{M} = \oplus$ – Les observations sont pilotées par une distribution de Bernoulli et deux gaussiennes

$$\begin{cases} \mathbf{X} | \mathbf{B} = \mathbf{m}, \boldsymbol{\theta}, \mathfrak{M} = \oplus & \sim \mathcal{N}(\mathbf{x}; \mathbf{m}_{\mathbf{m}}, \Gamma_{\mathbf{m}}^{-1}) \\ \mathbf{X} | \mathbf{B} = \mathbf{s}, \boldsymbol{\theta}, \mathfrak{M} = \oplus & \sim \mathcal{N}(\mathbf{x}; \mathbf{m}_{\mathbf{s}}, \Gamma_{\mathbf{s}}^{-1}) \\ \mathbf{B} | \boldsymbol{\theta}, \mathfrak{M} = \oplus & \sim \mathcal{B}(\mathbf{b}; p) \end{cases}$$

La densité pour la concentration, marginale vis-à-vis du statut, est alors un mélange de deux gaussiennes. Ces distributions font intervenir les paramètres inconnus $\boldsymbol{\theta} = [\mathbf{m}_{\mathbf{m}}, \mathbf{m}_{\mathbf{s}}, \Gamma_{\mathbf{m}}, \Gamma_{\mathbf{s}}, p]$ que l'on probabilise sous une a priori séparée $\pi_{\Theta}(\boldsymbol{\theta}) = \pi_{\mathbf{m}}(\mathbf{m}_{\mathbf{m}}, \Gamma_{\mathbf{m}}) \pi_{\mathbf{s}}(\mathbf{m}_{\mathbf{s}}, \Gamma_{\mathbf{s}}) \pi_P(p)$ avec

$$\begin{cases} (\mathbf{m}_{\mathbf{m}}, \Gamma_{\mathbf{m}}) | \mathfrak{M} = \oplus & \sim \mathcal{NW}(\mathbf{m}_{\mathbf{m}}, \Gamma_{\mathbf{m}}; \nu_{\mathbf{m}}, \eta_{\mathbf{m}}, \boldsymbol{\mu}_{\mathbf{m}}, \boldsymbol{\Lambda}_{\mathbf{m}}) \\ (\mathbf{m}_{\mathbf{s}}, \Gamma_{\mathbf{s}}) | \mathfrak{M} = \oplus & \sim \mathcal{NW}(\mathbf{m}_{\mathbf{s}}, \Gamma_{\mathbf{s}}; \nu_{\mathbf{s}}, \eta_{\mathbf{s}}, \boldsymbol{\mu}_{\mathbf{s}}, \boldsymbol{\Lambda}_{\mathbf{s}}) \\ p | \mathfrak{M} = \oplus & \sim \mathcal{P}(p; \alpha, \beta) \end{cases}$$

- $\mathfrak{M} = \ominus$ – Les observations sont pilotées par une distribution de Bernoulli et une (seule) gaussienne

$$\begin{cases} \mathbf{X} | \mathbf{B} = \mathbf{m}, \boldsymbol{\theta}, \mathfrak{M} = \ominus & \sim \mathcal{N}(\mathbf{x}; \mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}) \\ \mathbf{X} | \mathbf{B} = \mathbf{s}, \boldsymbol{\theta}, \mathfrak{M} = \ominus & \sim \mathcal{N}(\mathbf{x}; \mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}) \\ \mathbf{B} | \boldsymbol{\theta}, \mathfrak{M} = \ominus & \sim \mathcal{B}(\mathbf{b}; p) \end{cases}$$

Ici, \mathbf{X} et \mathbf{B} sont indépendants conditionnellement à \mathfrak{M} ce qui est bien conforme à l'idée que le jeu de protéines ne permet pas de discriminer. Ces distributions font intervenir les paramètres inconnus $\boldsymbol{\theta} = [\mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}, p]$ que l'on probabilise également sous une a priori séparée $\pi_{\Theta}(\boldsymbol{\theta}) = \pi_{\mathbf{c}}(\mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}) \pi_P(p)$ avec

$$\begin{cases} (\mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}) | \mathfrak{M} = \ominus & \sim \mathcal{NW}(\mathbf{m}_{\mathbf{c}}, \Gamma_{\mathbf{c}}; \nu_{\mathbf{c}}, \eta_{\mathbf{c}}, \boldsymbol{\mu}_{\mathbf{c}}, \boldsymbol{\Lambda}_{\mathbf{c}}) \\ p | \mathfrak{M} = \ominus & \sim \mathcal{P}(p; \alpha, \beta) \end{cases}$$

On probabilise aussi le modèle lui-même : $\Pr[\mathfrak{M} = \oplus]$ et $\Pr[\mathfrak{M} = \ominus]$ et on prendra 1/2 pour chacune.

¹La question suivante concerne la sélection des protéines les plus discriminantes mais ne sera pas détaillée ici. Les questions d'apprentissage en tant que tel et de classification sortent du cadre de ce travail.

Décision optimale et cote a posteriori

	$\mathfrak{m}^* = \oplus$	$\mathfrak{m}^* = \ominus$
$\widehat{\mathfrak{m}} = \oplus$	VD / $\mathcal{C}(\oplus, \oplus)$	FD / $\mathcal{C}(\oplus, \ominus)$
$\widehat{\mathfrak{m}} = \ominus$	FND / $\mathcal{C}(\ominus, \oplus)$	VND / $\mathcal{C}(\ominus, \ominus)$

TAB. 1 – Résultats de décision et coûts associés.

La suite du jeu consiste à construire un « décideur », *i.e.*, une fonction ψ de l'ensemble des observations vers l'ensemble des décisions. Pour cela on s'intéresse aux quatre situations possibles (voir Tab. 1) : deux décisions correctes (vrai discriminant et vrai non-discriminant) et deux erronées (faux discriminant et faux non-discriminant). On se alors un coût (indiqué dans Tab. 1) qui quantifie l'impact de chacune des situations. On s'intéresse au cas particulier : $\mathcal{C}(\oplus, \oplus) = \mathcal{C}(\ominus, \ominus) = 0$ pour les bonnes décisions et au cas symétrique $\mathcal{C}(\ominus, \oplus) = \mathcal{C}(\oplus, \ominus)$ pour les mauvaises. On définit alors le risque $\mathcal{R}(\psi)$ pour ψ en prenant le coût moyen :

$$\mathcal{R}(\psi) = \mathbb{E}[\mathcal{C}(\psi(\mathcal{X}, \mathbf{b}), \mathfrak{m})] \quad (1)$$

et il s'agit d'une *triple* moyenne, sur le modèle \mathfrak{m} , sur les paramètres θ et sur les observations $(\mathcal{X}, \mathbf{b})$. Il mesure globalement la performance de ψ et on choisit le décideur optimal ψ_{opt} comme son minimiseur. On montre qu'il repose sur la cote a posteriori, *i.e.*, le rapport des probabilités a posteriori des deux modèles :

$$\rho = \rho(\mathcal{X}, \mathbf{b}) = \Pr[\mathfrak{m} = \oplus | \mathcal{X}, \mathbf{b}] / \Pr[\mathfrak{m} = \ominus | \mathcal{X}, \mathbf{b}] \quad (2)$$

et qu'il choisit sa décision en comparant ρ au seuil $\sigma = [\mathcal{C}(\ominus, \oplus) - \mathcal{C}(\oplus, \oplus)] / [\mathcal{C}(\oplus, \ominus) - \mathcal{C}(\ominus, \ominus)]$

$$\psi(\mathcal{X}, \mathbf{b}) = \oplus \text{ si } \rho(\mathcal{X}, \mathbf{b}) > \sigma \quad \text{et} \quad \psi(\mathcal{X}, \mathbf{b}) = \ominus \text{ si } \rho(\mathcal{X}, \mathbf{b}) < \sigma$$

Ce seuil vaut 1 dans notre cas particulier, *i.e.*, l'hypothèse la plus probable a posteriori emporte la décision.

Ce résultat est très puissant car il assure la supériorité de ce décideur sur tous les autres : il n'existe pas de procédure, pas d'algorithme, pas de code informatique, pas de méthode... qui produise un meilleur résultat et cela appelle un commentaire. L'optimalité est obtenue au sens très particulier de (1) et il est construit dans un cadre fixé entre autres par les distributions choisies : une modification de ces distributions entraîne naturellement une modification de l'estimateur. De plus, ces distributions font apparaître des hyperparamètres qui peuvent être délicats à choisir et nous verrons également que l'idée des distributions non-informatives génère des difficultés.

Cote a posteriori et facteur de Bayes

Intéressons-nous à la cote a posteriori (2). En multipliant numérateur et dénominateur par la distribution pour $(\mathcal{X}, \mathbf{b})$ on fait apparaître la distribution jointe pour $(\mathfrak{m}, \mathcal{X}, \mathbf{b})$ que l'on refactorise en conditionnant par \mathfrak{m} :

$$\rho = \underbrace{\frac{\Pr[\mathfrak{m} = \oplus | \mathcal{X}, \mathbf{b}]}{\Pr[\mathfrak{m} = \ominus | \mathcal{X}, \mathbf{b}]}}_{\text{Cote a posteriori}} = \underbrace{\frac{\Pr[\mathfrak{m} = \oplus, \mathcal{X}, \mathbf{b}]}{\Pr[\mathfrak{m} = \ominus, \mathcal{X}, \mathbf{b}]}}_{\text{Facteur de Bayes}} = \underbrace{\frac{f_{\mathcal{X}, \mathbf{B} | \mathfrak{m}}(\mathcal{X}, \mathbf{b} | \mathfrak{m} = \oplus)}{\Pr[\mathfrak{m} = \oplus]}}_{\text{Cote a priori}} \cdot \underbrace{\frac{\Pr[\mathfrak{m} = \oplus]}{\Pr[\mathfrak{m} = \ominus]}}_{\text{Cote a priori}}. \quad (3)$$

On voit apparaître une quantité clé : le facteur de Bayes *i.e.*, le rapport des évidences des deux modèles. Ces évidences ont la forme d'une vraisemblance : elles décrivent la distribution des observations sachant la quantité d'intérêt et mesurent ainsi une adéquation observations – modèle. On les obtient par marginalisation

$$f_{\mathcal{X}, \mathbf{B} | \mathfrak{m}}(\mathcal{X}, \mathbf{b} | \mathfrak{m} = \odot) = \int_{\theta} f_{\mathcal{X}, \mathbf{B}, \Theta | \mathfrak{m}}(\mathcal{X}, \mathbf{b}, \theta | \mathfrak{m} = \odot) d\theta$$

avec $\odot = \ominus$ et $\odot = \oplus$. Cette opération peut s'avérer délicate dans certains cas et nécessiter des calculs numériques intensifs, par exemple par échantillonnage stochastique. Ici, les choses se déroulent plus aisément grâce aux choix pour les diverses lois. La quantité à marginaliser se décompose en une vraisemblance et une a priori :

$$f_{\mathcal{X}, \mathbf{B}, \Theta | \mathfrak{M}}(\mathcal{X}, \mathbf{b}, \theta | \mathfrak{M} = \odot) = f_{\mathcal{X}, \mathbf{B} | \Theta, \mathfrak{M}}(\mathcal{X}, \mathbf{b} | \theta, \mathfrak{M} = \odot) \pi_{\Theta | \mathfrak{M}}(\theta | \mathfrak{M} = \odot)$$

qui sont décrites ci-après. Pour $\mathfrak{M} = \oplus$, on voit apparaître deux gaussiennes (malades et sains) alors que pour $\mathfrak{M} = \ominus$ tous les individus sont regroupés sous une unique gaussienne.

- Pour ce qui est de la vraisemblance des paramètres attachée aux observations, on a

$$\begin{cases} f_{\mathcal{X}, \mathbf{B} | \Theta, \mathfrak{M}}(\mathcal{X}, \mathbf{b} | \theta, \mathfrak{M} = \oplus) = \left[\prod_{\mathcal{I}_m} \mathcal{N}(\mathbf{x}_n; \mathbf{m}_m, \Gamma_m) \right] \left[\prod_{\mathcal{I}_s} \mathcal{N}(\mathbf{x}_n; \mathbf{m}_s, \Gamma_s) \right] [p^{N_m}(1-p)^{N_s}] \\ f_{\mathcal{X}, \mathbf{B} | \Theta, \mathfrak{M}}(\mathcal{X}, \mathbf{b} | \theta, \mathfrak{M} = \ominus) = \left[\prod_{\mathcal{I}_m \mathcal{I}_s} \mathcal{N}(\mathbf{x}_n; \mathbf{m}_c, \Gamma_c) \right] [p^{N_m}(1-p)^{N_s}] \end{cases} \quad (4)$$

- Pour compléter le tableau, concernant les deux densités a priori pour les paramètres

$$\begin{cases} \pi_{\Theta | \mathfrak{M}}(\theta | \mathfrak{M} = \oplus) = \mathcal{NW}(\mathbf{m}_m, \Gamma_m; \nu_m, \eta_m, \boldsymbol{\mu}_m, \Lambda_m) \mathcal{NW}(\mathbf{m}_s, \Gamma_s; \nu_s, \eta_s, \boldsymbol{\mu}_s, \Lambda_s) \mathcal{P}(p; \alpha, \beta) \\ \pi_{\Theta | \mathfrak{M}}(\theta | \mathfrak{M} = \ominus) = \mathcal{NW}(\mathbf{m}_c, \Gamma_c; \nu_c, \eta_c, \boldsymbol{\mu}_c, \Lambda_c) \mathcal{P}(p; \alpha, \beta) \end{cases} \quad (5)$$

Effectuant les produits (4)-(5), on fait apparaître six facteurs ($\mathfrak{M} = \oplus$) ou quatre ($\mathfrak{M} = \ominus$). Un élément technique crucial est la possibilité d'intégrer en s'appuyant sur les propriétés des densités Gauss-Wishart et Bêta :

$$\begin{aligned} f_{\mathcal{X}, \mathbf{B} | \mathfrak{M}}(\mathcal{X}, \mathbf{b} | \mathfrak{M} = \oplus) &= \frac{K_{\mathcal{NW}}(\nu_m^{\text{pst}}, \eta_m^{\text{pst}}, \boldsymbol{\mu}_m^{\text{pst}}, \Lambda_m^{\text{pst}})}{K_{\mathcal{NW}}(\nu_m, \eta_m, \boldsymbol{\mu}_m, \Lambda_m)} \frac{K_{\mathcal{NW}}(\nu_s^{\text{pst}}, \eta_s^{\text{pst}}, \boldsymbol{\mu}_s^{\text{pst}}, \Lambda_s^{\text{pst}})}{K_{\mathcal{NW}}(\nu_s, \eta_s, \boldsymbol{\mu}_s, \Lambda_s)} \frac{K_{\mathcal{B}}(\alpha^{\text{pst}}, \beta^{\text{pst}})}{K_{\mathcal{B}}(\alpha, \beta)} \\ f_{\mathcal{X}, \mathbf{B} | \mathfrak{M}}(\mathcal{X}, \mathbf{b} | \mathfrak{M} = \ominus) &= \frac{K_{\mathcal{NW}}(\nu_c^{\text{pst}}, \eta_c^{\text{pst}}, \boldsymbol{\mu}_c^{\text{pst}}, \Lambda_c^{\text{pst}})}{K_{\mathcal{NW}}(\nu_c, \eta_c, \boldsymbol{\mu}_c, \Lambda_c)} \frac{K_{\mathcal{B}}(\alpha^{\text{pst}}, \beta^{\text{pst}})}{K_{\mathcal{B}}(\alpha, \beta)} \end{aligned}$$

qui font apparaître les constantes de normalisation pour les paramètres a posteriori et a priori. Pour avoir la cote a posteriori (3) il suffit de faire le ratio de ces rapports de constantes de normalisation. Nous avons là un résultat tout à fait remarquable puisqu'il fournit la cote a posteriori sans calcul numérique intensif et il suffit de calculer les paramètres a posteriori et les constantes de normalisation, comme donnés en annexes, Eq. (6)-(7). Ils font apparaître, *in fine* les moyennes et covariances empiriques (données Eq. (8)) qui interviennent dans les $\boldsymbol{\mu}_x^{\text{pst}}$ et Λ_x^{pst} .

Application au cas d'une unique protéine

Le résultat précédent est valide en dimension P quelconque et inclut le cas de concentrations corrélées. Cette section en présente une première illustration dans le cas $P = 1$. L'expression devient alors :

$$\rho \propto \frac{B_P(N_m/2, N_s/2)}{(1-\bar{p})^{PN_s/2} \bar{p}^{PN_m/2}} \left(1 + \frac{1-\bar{p}}{\bar{p}} \frac{\bar{\mathbf{R}}_s}{\bar{\mathbf{R}}_m} \right)^{N_m/2} \left(1 + \frac{\bar{p}}{1-\bar{p}} \frac{\bar{\mathbf{R}}_m}{\bar{\mathbf{R}}_s} \right)^{N_s/2} \left(1 + \bar{p}(1-\bar{p}) \frac{(\bar{\mathbf{x}}_m - \bar{\mathbf{x}}_s)^2}{\bar{p}\bar{\mathbf{R}}_m + (1-\bar{p})\bar{\mathbf{R}}_s} \right)^{N/2}$$

et elle fait naturellement apparaître les proportions, les moyennes et les variances empirique (voir annexe).

Remarque — *L'expression de ρ est donnée à un facteur multiplicatif près. Celui-ci est une fonction uniquement des hyperparamètres et peut-être difficile à manipuler lorsque l'on considère une a priori non-informative. Cela étant, il n'influence que la valeur du seuil dans la procédure de décision.*

Quoiqu'il en soit, le logarithme de la cote a posteriori en faveur de $\mathfrak{M} = \oplus$ (à un terme additif près) est présenté Fig. 1 comme fonction de la différence entre moyennes pour plusieurs variances et avec $N_m = 50$ et $N_s = 100$.

On observe un premier résultat attendu, conforme à l'intuition et à l'existant : pour des variances fixées, la cote augmente avec la différence des moyennes. En revanche, pour une différence de moyennes fixée, la cote n'augmente pas systématiquement lorsque la variance diminue ce qui peut sembler contre-intuitif. En réalité, lorsque la différence de moyenne est très faible (ou nulle à la limite), c'est la différence (ou le rapport) des variances qui peut rendre le jeu de protéines discriminant.

Cette analyse est préliminaire et les résultats sont en cours d'exploration plus avancée, notamment dans le cas multidimensionnel et corrélé.

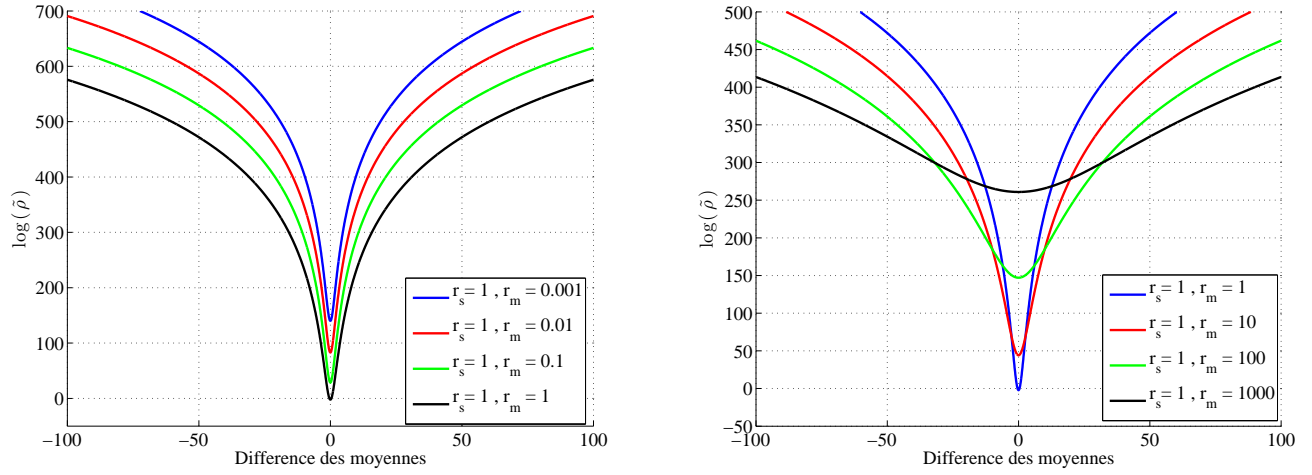


FIG. 1 – Évolution de la log-cote (à un terme additif près) avec la différence des moyennes. $N_m = 50$ et $N_s = 100$.

Annexe : trois résultats utiles

Densité Gauss-Wishart — Cette densité pour le couple $(\mathbf{m}, \mathbf{\Gamma})$ est paramétrée par quatre paramètres $\nu, \eta, \boldsymbol{\mu}, \mathbf{\Lambda}$:

$$\mathcal{NW}(\mathbf{m}, \mathbf{\Gamma}; \nu, \eta, \boldsymbol{\mu}, \mathbf{\Lambda}) = K_{\mathcal{NW}}^{-1} \det[\mathbf{\Gamma}]^{(\nu-P)/2} \exp \left[- \left(\text{tr}[\mathbf{\Gamma}\mathbf{\Lambda}^{-1}] + \eta(\mathbf{m} - \boldsymbol{\mu})^t \mathbf{\Gamma}(\mathbf{m} - \boldsymbol{\mu}) \right) / 2 \right] \quad (6)$$

et sa constante de normalisation s'écrit : $K_{\mathcal{NW}}(\nu, \eta, \boldsymbol{\mu}, \mathbf{\Lambda}) = (2\pi)^{P/2} 2^{\nu P/2} \eta^{-P/2} \det[\mathbf{\Lambda}]^{\nu/2} \Gamma_P(\nu/2)$ où Γ_P .

Paramètres a posteriori — Les paramètres a posteriori s'écrivent simplement

$$\nu_b^{\text{pst}} = N_b, \quad \eta_b^{\text{pst}} = N_b, \quad \boldsymbol{\mu}_b^{\text{pst}} = \bar{\mathbf{x}}_b, \quad (\mathbf{\Lambda}_b^{\text{pst}})^{-1} = N_b \bar{\mathbf{R}}_b \quad (7)$$

pour $\mathbf{b} = \mathbf{m}, \mathbf{b} = \mathbf{s}$ et $\mathbf{b} = \mathbf{c}$, dans le cas d'un a priori diffus (sans inclure d'information fiable sur les paramètres).

Moyennes et variances empiriques — On voit apparaître deux quantités attendues : moyenne et covariance empiriques de chaque groupe $\mathbf{b} = \mathbf{m}, \mathbf{b} = \mathbf{s}$ et $\mathbf{b} = \mathbf{c}$:

$$\bar{\mathbf{x}}_b = \frac{1}{N_b} \sum_{n \in \mathcal{I}_b} \mathbf{x}_n \quad \text{et} \quad \bar{\mathbf{R}}_b = \frac{1}{N_b} \sum_{n \in \mathcal{I}_b} (\mathbf{x}_n - \bar{\mathbf{x}}_b)(\mathbf{x}_n - \bar{\mathbf{x}}_b)^t. \quad (8)$$

On note $\bar{p} = N_m/N$ et $1 - \bar{p} = N_s/N$ les proportions observées d'individus malades et sains.

Références

- [1] C. P. Robert, *The Bayesian Choice. From decision-theoretic foundations to computational implementation*, Springer Texts in Statistics. Springer Verlag, New York, NY, USA, 2007.
- [2] P. Szacherski, J.-F. Giovannelli et P. Grangeat, « Joint Bayesian hierarchical inversion-classification and application in proteomics. », in *Proc. of the Int. Conf. on Stat. Signal Proc.*, Nice, France, juin 2011.
- [3] P. Szacherski, J.-F. Giovannelli, L. Gerfault, A. Giremus et P. Grangeat, « Robust MS serum sample classification in proteomics by the use of inverse problems », in *2012 IEEE International Workshop on Genomic Signal Processing and Statistics*, Washington, DC, USA, décembre 2012.