

Plans d'expérience pour l'estimation de la courbe de concentration et de l'AUC

Mohamad BELOUNI et Karim BENCHENNI

**Journée de rentrée des Doctorants du département
Probabilités-Statistique du LJK, 20 novembre 2012**

Laboratoire Jean Kuntzmann, Equipe IPS

Grenoble, France

**Ecole Doctorale Mathématiques, Sciences et Technologies
de l'Information, Informatique**

Université Joseph Fourier

- 1 **Définition du modèle et plans d'échantillonnages**
- 2 **Estimation de la fonction de concentration et de l'AUC**
- 3 **Perspectives**

Pour une fonction de concentration g définie dans un intervalle d'observation $[t_0, t_n]$, on considère le modèle :

$$X(t_i) = g(t_i) + \varepsilon(t_i) \quad \forall i = 0, 1, \dots, n,$$

où $\{t_0, t_1, \dots, t_n\}$ sont les instants d'observation de X dans l'intervalle $[t_0, t_n]$ et les erreurs $\varepsilon(t_i)$ forment un processus aléatoire autocorrélé de moyenne nulle $\mathbb{E}(\varepsilon(t)) = 0$ et de fonction d'autocovariance connue $\mathbb{E}(\varepsilon(t)\varepsilon(s)) = R(t, s)$.

On fixe les points $t_0 = 0$ et $t_n = 1$, et on s'intéresse au choix optimal des points intermédiaires $\{t_1, \dots, t_{n-1}\}$ pour l'estimation de l'intégrale (AUC) suivante :

$$I(g) = \int_0^1 g(t) dt.$$

On note l'estimateur linéaire de $I(g)$ par $L_n(g)$. Le but est de sélectionner les temps d'observations $\{t_1^*, \dots, t_{n-1}^*\}$ qui minimisent l'erreur quadratique moyenne

$$\text{EQM}(L_n(g)) = \mathbb{E}(L_n(g) - I(g))^2.$$

Plans d'échantillonnage réguliers de taille $n + 1$:

On suppose que $h(t)$ est une densité positive sur $[0, 1]$ avec une fonction de répartition strictement positive

$$H(t) = \int_0^t h(s) ds, \quad 0 \leq t \leq 1.$$

Le plan d'échantillonnage régulier défini par :

$$T_n = \left\{ t_{n,k} = H^{-1} \left(\frac{k}{n} \right), k = 0, \dots, n \right\},$$

comprend les bornes de l'intervalle $[0, 1]$, $t_{n,0} = 0$ et $t_{n,n} = 1$. Pour une densité uniforme sur $[0, 1]$, $h \equiv 1_{[0,1]}$, le plan d'échantillonnage régulier est périodique.

1 Définition du modèle et plans d'échantillonnages

2 Estimation de la fonction de concentration et de l'AUC

- Modèle de régression linéaire multiple
- Estimateur BLUE
- Estimateur simple et plan d'échantillonnage optimal

3 Perspectives

1 Définition du modèle et plans d'échantillonnages

2 Estimation de la fonction de concentration et de l'AUC

- **Modèle de régression linéaire multiple**
- Estimateur BLUE
- Estimateur simple et plan d'échantillonnage optimal

3 Perspectives

Modèle de régression linéaire multiple

$$X(t) = \beta' \mathbf{f}(t) + \varepsilon(t), \quad t \in [0, 1],$$

où le processus d'erreur $\varepsilon(t)$ est centré et d'autocovariance $\mathbb{E}(\varepsilon(t)\varepsilon(s)) = R(t, s)$,

$\beta' = (\beta_1, \dots, \beta_q)$ est un vecteur de paramètres inconnus et

$\mathbf{f}'(t) = (f_1(t), \dots, f_q(t))$ est un vecteur de fonctions de régression connues de la forme

$$f_j(t) = \int_0^1 R(t, s) \varphi_j(s) ds, \quad j = 1, \dots, q \quad t \in [0, 1] \quad (1)$$

où $\varphi_j, j = 1, \dots, q$, sont des fonctions continues sur $[0, 1]$.

- (H1) Sur la diagonale ($t = s$), on suppose que

$$R^{0,1}(t, t-) = \lim_{s \downarrow t} \partial R(t, s) / \partial s, \quad R^{0,1}(t, t+) = \lim_{s \uparrow t} \partial R(t, s) / \partial s,$$

existent et sont continues et

$$\alpha(t) = R^{0,1}(t, t-) - R^{0,1}(t, t+),$$

est supposée être strictement positive.

- (H2) En dehors de la diagonale ($t \neq s$) du carré unité, $R(t, s)$ est supposée avoir des dérivées partielles mixtes continues jusqu'à l'ordre deux .

1 Définition du modèle et plans d'échantillonnages

2 Estimation de la fonction de concentration et de l'AUC

- Modèle de régression linéaire multiple
- **Estimateur BLUE**
- Estimateur simple et plan d'échantillonnage optimal

3 Perspectives

Estimateur BLUE

On veut estimer la surface (AUC) sous la courbe de la fonction de concentration $g(t) = \beta' f(t)$ définie par

$$\text{AUC} = I(g) = \int_0^1 g(t) dt.$$

L'estimateur BLUE de β sur tout l'intervalle d'observation $[0, 1]$ est défini par

$$\hat{\beta} = \mathbf{S}^{-1} \int_0^1 X(t) \varphi(t) dt,$$

où $\varphi'(t) = (\varphi_1(t), \dots, \varphi_q(t))$, et la matrice \mathbf{S} de termes généraux

$$s_{ij} = \int_0^1 \int_0^1 \varphi_i(t) R(t, s) \varphi_j(s) dt ds, \quad i, j = 1, \dots, q$$

est supposée être inversible. On estime alors l'AUC ($I(g)$) par :

$$L(g) = \mathbf{z}' \hat{\beta},$$

où $\mathbf{z}' = (z_1, \dots, z_q)$; $z_j = \int_0^1 f_j(t) dt$, $j = 1, \dots, q$.

Si $X(t)$ est observé en $n + 1$ points d'un plan d'échantillonnage régulier $\mathbf{T}_n = \{t_{n,k}\}_0^n$ dans l'intervalle $[0, 1]$, alors l'estimateur BLUE $\hat{\beta}_n^{\text{blue}}$ de β du paramètre β défini par :

$$\hat{\beta}_n^{\text{blue}} = \mathbf{A}_{\mathbf{T}_n}^{-1} \boldsymbol{\eta},$$

où

$$\mathbf{A}_{\mathbf{T}_n} = (a_{rs})_{q \times q},$$

est supposée être inversible et de termes généraux

$$a_{rs} = \sum_{i,j=0}^n f_r(t_i) \mathbf{R}_{\mathbf{T}_n}^{-1}(t_i, t_j) f_s(t_j) \quad r, s = 1, \dots, q \text{ et}$$

$$\boldsymbol{\eta} = (\eta_r)_{q \times 1},$$

où

$$\eta_r = \sum_{i,j=0}^n f_r(t_i) \mathbf{R}_{\mathbf{T}_n}^{-1}(t_i, t_j) X(t_j), \quad r = 1, \dots, q$$

et $\mathbf{R}_{\mathbf{T}_n} = (\mathbf{R}_{\mathbf{T}_n}(t_{n,k}, t_{n,j}))_{(n+1) \times (n+1)}$ est supposée être inversible pour tout n .

alors on introduit l'estimateur de l'AUC ($l(g)$) par :

$$L_n^{\text{blue}}(g) = \mathbf{z}' \hat{\beta}_n^{\text{blue}} = \mathbf{z}' \mathbf{A}_{T_n}^{-1} \eta.$$

Inconvenients :

Les estimateurs BLUE, $\hat{\beta}_n^{\text{blue}}$, $L_n^{\text{blue}}(g)$ sont **instables** car ils utilisent l'inverse de la matrice d'autocovariance.

Remède :

On construit deux estimateurs $\hat{\beta}_n^{\text{trap}}(h)$ et $L_n^{\text{trap}}(g)$ qui sont **simples et plus stables** car ils ne dépendent pas de la matrice d'autocovariance et qui utilisent seulement les observations aux points du plan d'échantillonnage régulier.

1 Définition du modèle et plans d'échantillonnages

2 Estimation de la fonction de concentration et de l'AUC

- Modèle de régression linéaire multiple
- Estimateur BLUE
- **Estimateur simple et plan d'échantillonnage optimal**

3 Perspectives

Estimateur simple et plan d'échantillonnage optimal

Si $X(t)$ est observé en $n + 1$ points d'un plan d'échantillonnage régulier $\mathbf{T}_n = \{t_{n,k}\}_0^n$ dans l'intervalle $[0, 1]$, alors on construit l'estimateur simple de β par :

$$\hat{\beta}_n^{\text{trap}}(h) = \mathbf{S}_n^{*-1} \mathbf{L}_{n,q}(X),$$

où

$$\mathbf{L}_{n,q}(X) = \frac{1}{2n} \sum_{k=0}^{n-1} \left(\left(\frac{\varphi X}{h} \right)(t_k) + \left(\frac{\varphi X}{h} \right)(t_{k+1}) \right),$$

est un vecteur colonne

$$\mathbf{S}_n^* = (\mathbf{s}_{ij}^*)_{q \times q},$$

est supposée être inversible,

$$\mathbf{s}_{ij}^* = \frac{1}{2n} \sum_{k=0}^{n-1} \left(\left(\frac{\varphi_i}{h} f_j \right)(t_k) + \left(\frac{\varphi_i}{h} f_j \right)(t_{k+1}) \right), \quad i, j = 1, \dots, q$$

où $f_j(t), j = 1, \dots, q$ est définie par (1)

alors on introduit l'estimateur de l'AUC ($I(g)$) par :

$$L_n^{\text{trap}}(g) = \mathbf{z}' \hat{\beta}_n^{\text{trap}}(h) = \mathbf{z}' \mathbf{S}_n^{*-1} \mathbf{L}_{n,q}(X),$$

où $\mathbf{z}' = (z_1, \dots, z_q)$; $z_j = \int_0^1 f_j(t) dt$, $j = 1, \dots, q$.

Théorème

Si les hypothèses (H1) et (H2) sont vérifiées et si on suppose que $\varphi_i/h, i = 1, \dots, q$ est deux fois continûment différentiable, alors l'estimateur $\hat{\beta}_n^{trap}(h)$ vérifie

$$\lim_{n \rightarrow \infty} n^2 \left(EQM(\hat{\beta}_n^{trap}(h)) - trace(\mathbf{S}^{-1}) \right) = \frac{1}{12} \int_0^1 \alpha(t) \frac{\varphi'(t) \mathbf{S}^{-2} \varphi(t)}{h^2(t)} dt.$$

Par ailleurs, l'estimateur $\hat{\beta}_n^{trap}(h)$ avec un plan d'échantillonnage \mathbf{T}_n^* engendré par la densité

$$h^*(t) = \left\{ \alpha(t) \varphi'(t) \mathbf{S}^{-2} \varphi(t) \right\}^{1/3} / \int_0^1 \left\{ \alpha(u) \varphi'(u) \mathbf{S}^{-2} \varphi(u) \right\}^{1/3} du,$$

est asymptotiquement optimal.

Corollaire

L'estimateur $\hat{\beta}_n^{trap}(h^*)$, construit à partir de \mathbf{T}_n^* et engendré par la densité $h^*(t)$, vérifie :

$$\lim_{n \rightarrow \infty} n^2 \left(EQM(\hat{\beta}_n^{trap}(h^*)) - trace(\mathbf{S}^{-1}) \right) = \frac{1}{12} \left(\int_0^1 \left\{ \alpha(t) \varphi'(t) \mathbf{S}^{-2} \varphi(t) \right\}^{1/3} dt \right)^3$$

Théorème

Sous les hypothèses (H1) et (H2) et si on suppose que φ_i/h , $i = 1, \dots, q$ est deux fois continûment différentiable, alors l'estimateur $L_n^{trap}(g)$ de l'AUC, vérifie

$$\lim_{n \rightarrow \infty} n^2 \left(EQM(L_n^{trap}(g)) - \mathbf{z}' \mathbf{S}^{-1} \mathbf{z} \right) = \frac{1}{12} \int_0^1 \alpha(t) \frac{\varphi'(t) (\mathbf{S}^{-1} \mathbf{A} \mathbf{S}^{-1}) \varphi(t)}{h^2(t)} dt,$$

où $\mathbf{A} = \mathbf{z} \mathbf{z}'$.

Par ailleurs, l'estimateur $L_n^{trap}(g)$ avec un plan d'échantillonnage \mathbf{T}_n^* engendré par la densité $h^*(t)$

$$h^*(t) = \left\{ \alpha(t) \varphi'(t) (\mathbf{S}^{-1} \mathbf{A} \mathbf{S}^{-1}) \varphi(t) \right\}^{1/3} / \int_0^1 \left\{ \alpha(u) \varphi'(u) (\mathbf{S}^{-1} \mathbf{A} \mathbf{S}^{-1}) \varphi(u) \right\}^{1/3} du.$$

est asymptotiquement optimal.

Corollaire

L'estimateur $L_n^{\text{trap}}(g)$, construit à partir de \mathbf{T}_n^* et engendré par la densité $h^*(t)$, vérifie :

$$\lim_{n \rightarrow \infty} n^2 (EQM(L_n^{\text{trap}}(g)) - EQM(L(g))) = \frac{1}{12} \left(\int_0^1 \left\{ \alpha(t) \varphi'(t) (\mathbf{S}^{-1} \mathbf{A} \mathbf{S}^{-1}) \varphi(t) \right\}^{1/3} dt \right)^3$$

où $\mathbf{A} = \mathbf{z} \mathbf{z}'$.

Cas particulier : $q = 1$

On considère d'abord le modèle linéaire simple ($q = 1$)

$$X(t) = \beta f(t) + \varepsilon(t), \quad t \in [0, 1]$$

où $f(t)$ est de la la forme générale suivante :

$$f(t) = \int_0^1 R(t, s)\varphi(s) ds + \sum_{l=1}^L b_l R(t, a_l),$$

où φ est une fonction continue connue sur $[0, 1]$, les constantes $b_l, l = 1, \dots, L$ sont connues et non nulles et les $a_l, l = 1, \dots, L$ sont des points connus dans $[0, 1]$.

Si $X(t)$ est observé en $n + 1$ points d'un plan d'échantillonnage régulier $\mathbf{T}_n = \{t_{n,k}\}_0^n$ engendré par une densité positive $h(t)$, alors on construit l'estimateur simple de β est défini par :

$$\hat{\beta}_n^{\text{trap}} = \frac{L_n(X)}{L_n(f)},$$

où

$$L_n(X) = \frac{1}{2n} \sum_{k=0}^{n-1} \left(\left(\frac{\varphi}{h} X \right)(t_{n,k}) + \left(\frac{\varphi}{h} X \right)(t_{n,k+1}) \right) + \sum_{l=1}^L b_l X(a_l),$$

$$L_n(f) = \frac{1}{2n} \sum_{k=0}^{n-1} \left(\left(\frac{\varphi}{h} f \right)(t_{n,k}) + \left(\frac{\varphi}{h} f \right)(t_{n,k+1}) \right) + \sum_{j=1}^L b_j f(a_j).$$

$\text{Var}(\hat{\beta}_n^{\text{trap}}) \longrightarrow \text{Var}(\hat{\beta}) = s_L^{-2}$ quand $n \rightarrow \infty$, où

$$s_L^2 \triangleq \int_0^1 \int_0^1 \varphi(t) R(t, s) \varphi(s) dt ds + 2 \sum_{l=1}^L b_l \int_0^1 \varphi(s) R(s, a_l) ds + \sum_{l=1}^L \sum_{j=1}^L b_l R(a_l, a_j) b_j.$$

Théorème

Sous les hypothèses (H1) et (H2) et si on suppose que φ/h , $i = 1, \dots, q$ est deux fois continûment différentiable, alors l'estimateur $\hat{\beta}_n^{trap}$, vérifie

$$\lim_{n \rightarrow \infty} n^2 \left\{ \text{Var}(\hat{\beta}_n^{trap}) - s_L^{-2} \right\} = \frac{s_L^{-4}}{12} \int_0^1 \alpha(t) \frac{\varphi^2(t)}{h^2(t)} dt,$$

Par ailleurs, l'estimateur $\hat{\beta}_n^{trap}$ avec un plan d'échantillonnage \mathbf{T}_n^* engendré par la densité $h^*(t)$

$$h^*(t) = \left\{ \alpha(t) \varphi^2(t) \right\}^{1/3} / \int_0^1 \left\{ \alpha(u) \varphi^2(u) \right\}^{1/3} du,$$

est asymptotiquement optimal.

Références :

[1] SACKS, J. AND YLVISAKER, D (1966). Designs for regression problems with correlated errors, *Ann. Math*, 37, 66-89.

[2] CAMBANIS, S. AND SU, Y (1993). Sampling designs for regression coefficient estimation with correlated errors, *Annals of the Institute of Statistical Mathematics*, 46,707-722.

Exemple

On considère le modèle de régression simple

$$X(t) = \beta t + \varepsilon(t), \quad t \in [0, 1]$$

où $\varepsilon(t)$ est un processus d'erreur de type "Ornstein - Uhlenbeck" (Gauss-Markov) et sa fonction d'autocovariance $R(t, s) = \sigma^2 e^{-\lambda|t-s|}$. On montre que $f(t) = t$ peut s'écrire sous la forme suivante

$$f(t) = \int_0^1 R(t, s) \varphi(s) ds + \sum_{l=1}^L b_l R(t, a_l),$$

avec

$$\varphi(s) = \lambda s / (2\sigma^2), \quad L = 2, \quad a_1 = 0, \quad a_2 = 1, \quad b_1 = -1 / (2\lambda\sigma^2),$$

$$b_2 = (\lambda + 1) / (2\lambda\sigma^2).$$

La fonction de sauts (H1) est :

$$\alpha(t) = R^{0,1}(t, t-) - R^{0,1}(t, t+) = 2\lambda\sigma^2 > 0$$

La densité du plan d'échantillonnage (asymptotiquement) optimal $h^*(t)$ est

$$h^*(t) = (5/3)t^{2/3}$$

et les points optimaux sont

$$t_{n,k}^* = \{k/n\}^{3/5}, k = 0, \dots, n.$$

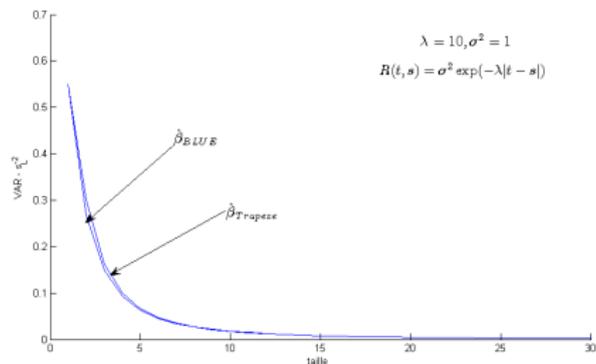
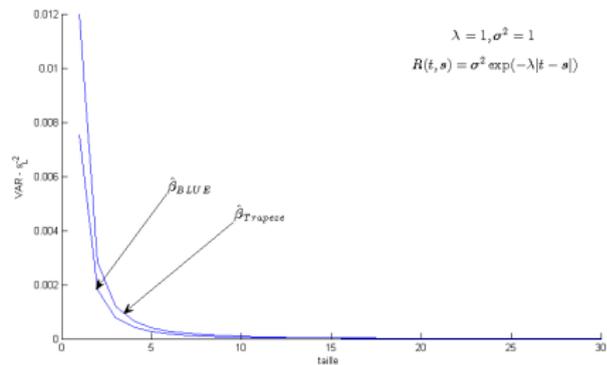


FIG.: $\left(\text{Var}(\hat{\beta}_n) - \text{Var}(\hat{\beta})\right)$ contre la taille de l'échantillon

1 Définition du modèle et plans d'échantillonnages

2 Estimation de la fonction de concentration et de l'AUC

- Modèle de régression linéaire multiple
- Estimateur BLUE
- Estimateur simple et plan d'échantillonnage optimal

3 Perspectives

- 1 Plans d'échantillonnage optimaux pour l'estimation non-paramétrique de la fonction de concentration et de l'aire sous la courbe de concentration (AUC)

En cours.

- 2 Taille fixée et petite : application à la pharmacocinétique (estimation + algorithmes d'optimalité)

En cours.