

PROCÉDURE DE TEST POUR L'AUTO-CALIBRATION BASÉE SUR LES COURBES DE LORENZ ET DE CONCENTRATION

Huyghe Julie¹ & Julien Trufin² & Thomas Verdebout³

¹ *Department of Mathematics Université Libre de Bruxelles (ULB) Brussels, Belgium - julie.huyghe@ulb.be*

² *Department of Mathematics Université Libre de Bruxelles (ULB) Brussels, Belgium - julien.trufin@ulb.be*

² *Department of Mathematics Université Libre de Bruxelles (ULB) Brussels, Belgium - thomas.verdebout@ulb.be*

1 Courbes de Lorentz, de concentration et auto-calibration.

Nous considérons dans ce travail un contexte de régression classique dans lequel on a (i) une variable aléatoire **réponse** Y à valeurs réelles et (ii) un ensemble de **caractéristiques** X_1, \dots, X_p regroupées dans un vecteur aléatoire \mathbf{X} à valeurs dans \mathbb{R}^p . La structure de dépendance à l'intérieur du vecteur aléatoire (Y, X_1, \dots, X_p) est ici exploitée pour extraire l'information contenue dans \mathbf{X} à propos de Y . Soit

$$\mu(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

l'espérance conditionnelle de Y sachant \mathbf{X} et soit

$$\begin{aligned} \hat{m} : \mathbb{R}^p &\rightarrow \mathbb{R} \\ \mathbf{x} &\rightarrow \hat{m}(\mathbf{x}) \end{aligned}$$

un estimateur de $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$. On peut associer à cet estimateur \hat{m} , deux courbes qui seront importantes ici:

- la **courbe de concentration** de $\mu(\mathbf{X})$ par rapport à $\hat{m}(\mathbf{X})$ définie par

$$\alpha \mapsto \text{CC}[\mu(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] = \frac{\mathbb{E} [\mu(\mathbf{X}) I[\hat{m}(\mathbf{X}) \leq F_{\hat{m}}^{-1}(\alpha)]]}{\mathbb{E}[\mu(\mathbf{X})]}, \quad \alpha \in (0, 1),$$

où $F_{\hat{m}}^{-1}$ est la fonction quantile associée à $\hat{m}(\mathbf{X})$ et

- la **courbe de Lorenz** associée à $\hat{m}(\mathbf{X})$ définie par

$$\begin{aligned} \alpha \mapsto \text{LC}[\hat{m}(\mathbf{X}); \alpha] &= \text{CC}[\hat{m}(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] \\ &= \frac{\mathbb{E} [\hat{m}(\mathbf{X}) I[\hat{m}(\mathbf{X}) \leq F_{\hat{m}}^{-1}(\alpha)]]}{\mathbb{E}[\hat{m}(\mathbf{X})]}, \quad \alpha \in (0, 1). \end{aligned}$$

En assurance, L'estimateur \hat{m} de la fonction de régression a pour objectif d'être utilisée comme un prime et il est donc naturel que la somme des primes correspondent le plus possible a la somme des pertes réelles. Ceci conduit naturellement au concept d'autocalibration. L'estimateur \hat{m} est dit auto-calibré si et seulement si

$$E[Y|\hat{m}(\mathbf{X}) = m] = m$$

pour tout $m \in \mathbb{R}$. Comme nous allons le décrire dans la prochaine Section, les courbes de concentration et de Lorenz vont être utilisée pour tester l'auto-calibration. Plus précisément, nous fournissons ci-dessous un test fondé sur la comparaison de fonctions aléatoires.

1.1 Procédure de test pour l'auto-calibration

Nous avons que

$$CC[\mu(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] = LC[\hat{m}(\mathbf{X}); \alpha] \text{ pour tout } \alpha$$

si et seulement si la version non-biaisée de \hat{m} donnée par $\hat{m}_{\text{unbiased}}(\mathbf{X}) := \frac{E[Y]}{E[\hat{m}(\mathbf{X})]} \hat{m}(\mathbf{X})$ est auto-calibrée. Nous utilisons ci-dessous cette caractérisation en comparant des version empirique des courbes de concentration et de Lorenz. Cette caractérisation de l'auto-calibration via les courbes de performance nous permet de proposer une procédure de test pour l'auto-calibration à partir de n copies i.i.d. $(Y_1, \mathbf{X}_1), \dots, (Y_n, \mathbf{X}_n)$ de (Y, \mathbf{X}) . Plus précisément, nous souhaitons tester l'hypothèse nulle :

$$\mathcal{H}_0 : CC[\mu(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] = LC[\hat{m}(\mathbf{X}); \alpha] \text{ pour tout } \alpha \in (0, 1)$$

contre l'hypothèse alternative :

$$\mathcal{H}_1 : CC[\mu(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] \neq LC[\hat{m}(\mathbf{X}); \alpha] \text{ pour un } \alpha \in (0, 1).$$

La procédure de test est donc naturellement basée sur la différence entre les versions empiriques des courbes de Lorenz et de concentration, définies comme suit :

$$\widehat{CC}[\mu(\mathbf{X}), \hat{m}(\mathbf{X}); \alpha] = \frac{1}{n\bar{Y}} \sum_{i=1}^n Y_i I [\hat{m}(\mathbf{X}_i) \leq F_{\hat{m}}^{-1}(\alpha)], \quad \alpha \in (0, 1),$$

et

$$\widehat{LC}[\hat{m}(\mathbf{X}); \alpha] = \frac{1}{n\bar{m}} \sum_{i=1}^n \hat{m}(\mathbf{X}_i) I [\hat{m}(\mathbf{X}_i) \leq F_{\hat{m}}^{-1}(\alpha)], \quad \alpha \in (0, 1).$$

Plus précisément, l'hypothèse nulle est rejetée pour de grandes valeurs de la statistique de test suivante :

$$\mathcal{T} = \sup_{\alpha \in (0, 1)} |T_n(\alpha)|,$$

où

$$\begin{aligned} T_n(\alpha) &= \sqrt{n} \left(\widehat{\text{CC}}[\mu(\mathbf{X}), \widehat{m}(\mathbf{X}); \alpha] - \widehat{\text{LC}}[\widehat{m}(\mathbf{X}); \alpha] \right) \\ &= n^{-1/2} \sum_{i=1}^n \left(\frac{Y_i}{\bar{Y}} - \frac{\widehat{m}(\mathbf{X}_i)}{\bar{m}} \right) I[\widehat{m}(\mathbf{X}_i) \leq F_{\widehat{m}}^{-1}(\alpha)] \end{aligned}$$

On peut ensuite dériver le comportement asymptotique du processus $T_n(\alpha)$. En supposant que $(Y_i, \widehat{m}(\mathbf{X}_i))$, $i = 1, 2, \dots, n$, sont tels que : $E[Y_i] \neq 0$, $E[\widehat{m}(\mathbf{X}_i)] \neq 0$, $E[\widehat{m}^2(\mathbf{X}_i)] < \infty$ et $E[Y_i^2] < \infty$. Alors sous l'hypothèse nulle, $T_n(\alpha)$ converge vers un processus gaussien de moyenne zero et de covariance

$$\mathcal{C}(\alpha_1, \alpha_2) = \mathbf{v}'(\alpha_1, \alpha_2) \boldsymbol{\Sigma}(\alpha_1, \alpha_2) \mathbf{v}(\alpha_1, \alpha_2).$$

avec (en supposant que toutes ces quantités existent et sont finies) :

$$Z_i(\alpha) = Y_i I[\widehat{m}(\mathbf{X}_i) \leq F_{\widehat{m}}^{-1}(\alpha)] \text{ and } W_i(\alpha) = \widehat{m}(\mathbf{X}_i) I[\widehat{m}(\mathbf{X}_i) \leq F_{\widehat{m}}^{-1}(\alpha)].$$

$\boldsymbol{\Sigma}(\alpha_1, \alpha_2)$ la matrice de covariance du vecteur aléatoire

$$(\widehat{m}(\mathbf{X}_i), Y_i, Z_i(\alpha_1), W_i(\alpha_1), Z_i(\alpha_2), W_i(\alpha_2))',$$

pour $(\alpha_1, \alpha_2) \in (0, 1)^2$, et enfin

$$\mathbf{v}(\alpha_1, \alpha_2) = \begin{pmatrix} -\frac{e_w(\alpha_1)}{e_m^2} & -\frac{e_z(\alpha_1)}{e_y^2} & e_y^{-1} & e_m^{-1} & 0 & 0 \\ -\frac{e_w(\alpha_2)}{e_m^2} & -\frac{e_z(\alpha_2)}{e_y^2} & 0 & 0 & e_y^{-1} & e_m^{-1} \end{pmatrix}',$$

où

$$e_y = E[Y_i], \quad e_m = E[\widehat{m}(\mathbf{X}_i)], \quad e_z(\alpha) = E[Z_i(\alpha)] \text{ and } e_w(\alpha) = E[W_i(\alpha)].$$

La procédure de test rejette donc l'hypothèse nulle si $\mathcal{T} > c_\beta$ où β est un niveau de confiance fixé tel que $P[\sup_{\alpha \in (0,1)} |T_n(\alpha)| > c_\beta] = \beta$ sous l'hypothèse nulle. On ne peut pas calculer c_β car la distribution de $(Y_i, \widehat{m}(\mathbf{X}_i))$ reste inconnue. Nous utilisons néanmoins ce résultat pour fournir une procédure bootstrap.

Bibliographie

Yitzhaki, S., Schechtman, E. (2013). The Gini Methodology: A Primer on Statistical Methodology. Springer.

Krüger, F., Ziegel, J.F. (2021). Generic conditions for forecast dominance. *Journal of Business & Economic Statistics* 39, 972-983.

Frees, E.W., Meyers, G., Cummings, A.D. (2011). Summarizing insurance scores using a Gini index. *Journal of the American Statistical Association* 106, 1085-1098.

Denuit, M., Sznajder, D., Trufin, J. (2019). Model selection based on Lorenz and Concentration curves, Gini indices and convex order. *Insurance: Mathematics and Economics* 89, 128-139.