

Calage à poids bornés : Que fait-on au juste ?

Jean-Claude DEVILLE (ex-Crest/Ensai)

Mise en place

Pour toute unité k de la population x_k est un vecteur de p variables auxiliaires présentes dans les données d'un échantillon de taille n .

Le total dans la population des x_k est connu et noté t .

On dispose d'un estimateur initial \hat{t} (en général Horvitz-Thompson, mais ça n'a rien d'obligatoire) linéaire utilisant des poids d_k .

X est la matrice $p \times n$ des $x_k d_k$. Elle est supposée de plein rang et, pour simplifier les notations, on écrira toujours x_k pour $x_k d_k$.

$H = \text{Ker}(X)$ est noyau de dimension $n-p$ de X .

On a donc $\hat{t} = X\mathbf{1}$ ($\mathbf{1}$ vecteur dont les n coordonnées valent 1)

Comment on fait ? On cherche :

-de nouveaux poids g_k compris entre deux valeurs m et M

- obtenir l'égalité $t = \sum_S x_k g_k = Xg$ (calage).

Calage classique : on minimise une 'distance' additive

$\sum_S G(g_k, 1)$ où G doit avoir de bonnes propriétés (convexité stricte et dérivabilité).

La situation est simple : soit le problème n'admet pas de solution, soit admet une solution unique donnée par la résolution des équations de calage $t = \sum_S x_k F(q_k x'_k \lambda)$ où λ est un vecteur de multiplicateurs de Lagrange et F une fonction définie sur \mathbb{R} déduite de la distance , croissante et 'régulière' (dérivable, de limites m et M les limites étant atteintes soit asymptotiquement soit à distance finie).

Remarque : les q_k furent introduit pour justifier l'estimateur ratio....avant qu'on s'aperçoive que c'est une variante de calage généralisé !!!

Si la distance est définie sur \mathbb{R} tout entier on a toujours une solution et il en va naturellement de même quand m et M sont suffisamment proches de $\pm\infty$. A l'inverse si $M-m$ est faible on n'aura généralement pas de solution. Si $M-m=0$, la seule solution possible est un estimateur par ratio !

Une pratique courante : chercher par tâtonnements des calages où m et M seront respectivement le plus grand et le plus petit possible.

On constate les choses suivantes :

- **M et m ne dépendent pas de la distance utilisée**
- **La plupart des poids convergent vers M ou m**

Pourquoi ?

En fait, on a un peu changé le problème, et on étudie la position relative de la variété linéaire $Xg=t$ (de dimension $n-p$) avec le n -cube centré en $c=(M+m)\mathbf{1}/2$ et de côté $D=M-m$.

Supposons c fixé pour commencer ; 3 cas possibles :

- L'intersection est vide (pour D assez petit).
- Elle contient au moins un point intérieur au cube (pour D assez grand) et elle est de dimension $n-p$.
- D est la valeur limite où l'intersection est un polytope convexe contenu dans une face de la frontière du cube (de dimension au plus $n-1$). Sa dimension peut varier de 0 (point unique) à $n-p$.

La procédure de tâtonnement consiste à rechercher un couple (c, D) et un vecteur 'calant' g tels que $\max_k (|g_k - c|)$ soit minimum.

Il est plus parlant de regarder ce qui se passe dans \mathbb{R}^p espace image de X . On étudie donc la position relative de t et de $K_{c,D}$ image par X du cube de centre c et de coté D .

C'est un polytope de dimension p , fermé convexe et symétrique de centre Xc . Il y a calage quand t est intérieur à $K_{c,D}$.

A c fixé, c'est toujours possible pour $D \geq D_c$:

-Si $<$ on minimise la distance (calage 'habituel')

-si = on obtient un point $t_g = Xg$ de la frontière de K_{c,D_c} résultat du tâtonnement. t_g est dans une $p-1$ (ou moins !) face de K_{c,D_c} qui est l'image d'une $p-1$ face du n -cube.

C'est aussi la limite de $F(q_k x'_k \lambda)$ pour λ infini.

On remarque maintenant que Xc (le centre) est autorisé à varier dans la droite qui contient 0 et $\hat{t} = X\mathbf{1}$. La minimisation de D_c achève le tâtonnement.

Il en résulte que $n-p+2$ des g_k valent m ou M ('équations' d'une $p-2$ face), les $p-2$ restant servant à définir notre point dans sa $p-2$ -face.

Comme annoncé par la pratique, ce résultat explique le groupement de presque tous les poids aux deux bornes, indépendamment de la distance utilisée pour faire les calculs !

La recherche de c est en fait l'estimation d'un paramètre supplémentaire qui amène notre point au bord de sa $p-1$ -face, sur une $p-2$ face avec un g_k de plus à la borne.

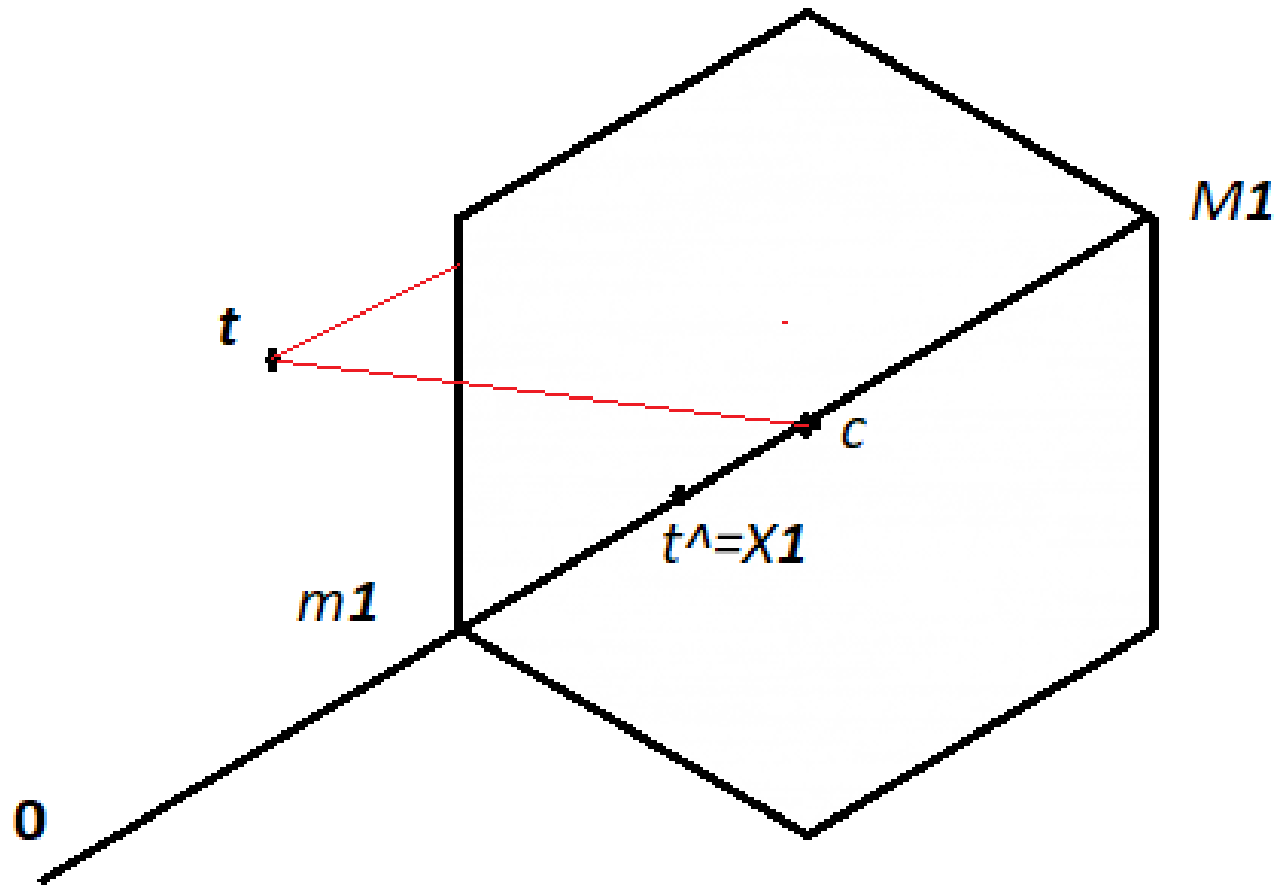


Illustration du bazar

Quelques commentaires

Le cas où l'intersection ne se réduit pas à un point est lié à une particularité des données (colinéarités pour des familles de moins de p vecteurs x_k). Je m'en tire par une construction qui les élimine mais c'est compliqué bien que ça ne change pas la nature des choses.

On a en fait estimé 2 paramètres de plus m et M . Quel sont les conséquences sur les justifications (asymptotiques !) du calage ?

Sans doute rien de dramatique, mais va t'en savoir...

On minimise une distance de type l_∞ . Approximable par l_p pour p grand ?

On pourrait généraliser à des situations plus complexes , par exemple avec des m_k et M_k (robustesse, domaines...).

Où les choses se compliquent un peu :Où on en est ??

1- Si t intérieur à K $t=Xg$ avec $g_k=F(q_k x'_k \lambda)$ et $\lambda \in \mathbb{R}^p$

2-Si t à la frontière de K presque tout les poids valent m ou M et correspondent à un $\lambda = \pm \infty$ (pour au moins une composante).

En fait les g admissibles constituent une variété G de dimension p telle que $G+H$ recouvre le cube C_n . Tout point du cube se projette le long de H en un point unique de G .

3- S_K = ensemble des sommets de K est l'image par X d'une partie des sommets de C_n , qui définissent le 'contour apparent'. Les autres se 'projettent' à l'intérieur de K . (ex :un 3-cube)

Cela vaut pour le calage 'métrique' où X détermine, avec la métrique, la variété G .

Le contour apparent, lui, ne dépend que de $\text{Ker}(X)$!!

Est-ce encore vrai dans le cas du calage généralisé , où G est de la forme $g_k = F(z'_k \lambda)$ ($\lambda \in \mathbb{R}^p$) avec des z_k ‘quelconques’ ?

NON !

Un exemple minimaliste :

$$X = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \text{ et } Z = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

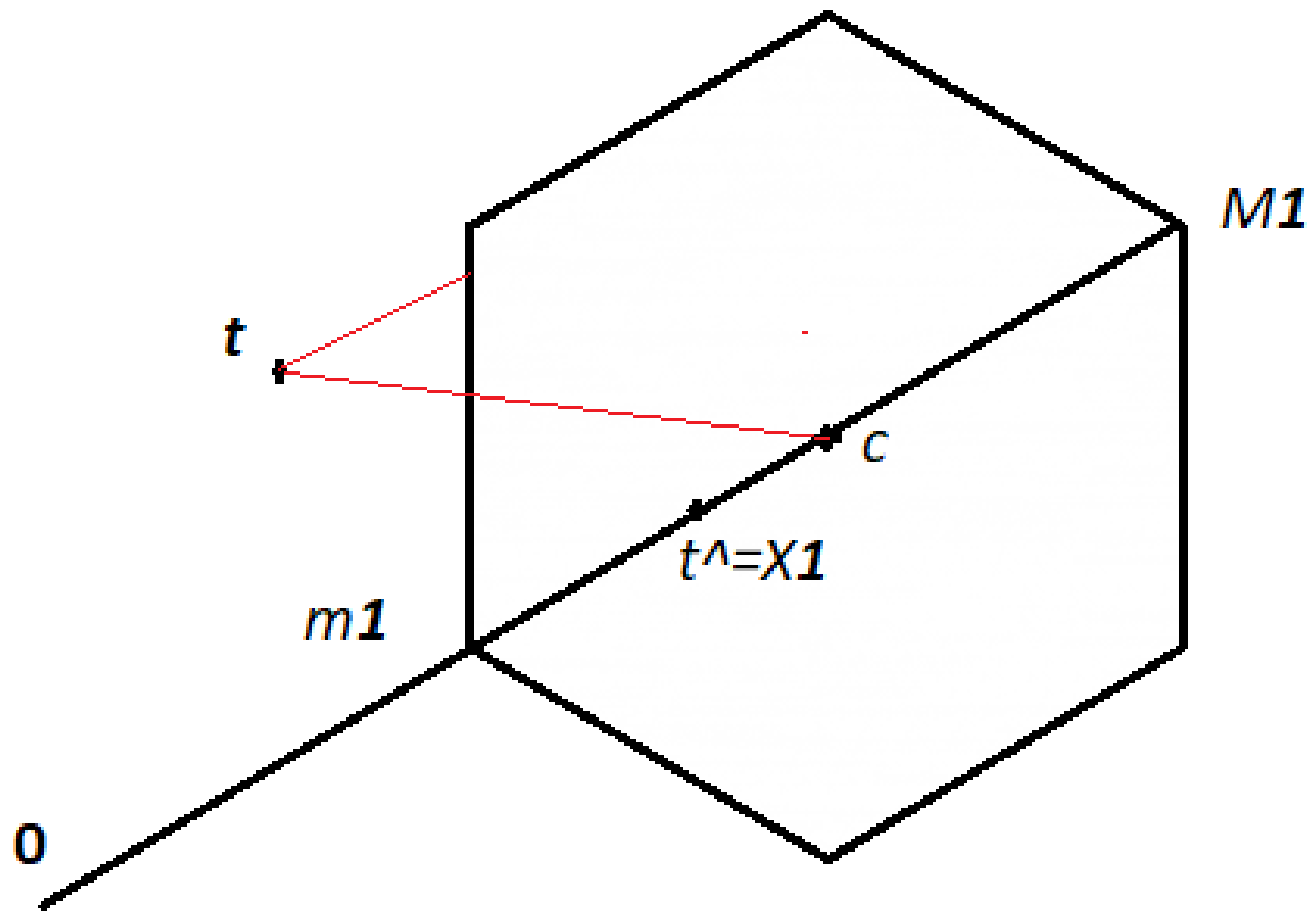
$$\text{Ker}(X) = (1 \ 1 \ -1)'$$

c'est du calage ordinaire et tout va bien , avec n'importe quelle fonction $F!$

Si on fait $X = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{pmatrix}$ c'est la catastrophe !

$$\text{Ker}(X) = (1 \ -1 \ -1)'$$

Les sommets $(1 \ 1 \ 0)$ et $(0 \ 0 \ 1)$ du cube ne sont pas atteints...et on trouve des cas où le calage a deux solutions distinctes.



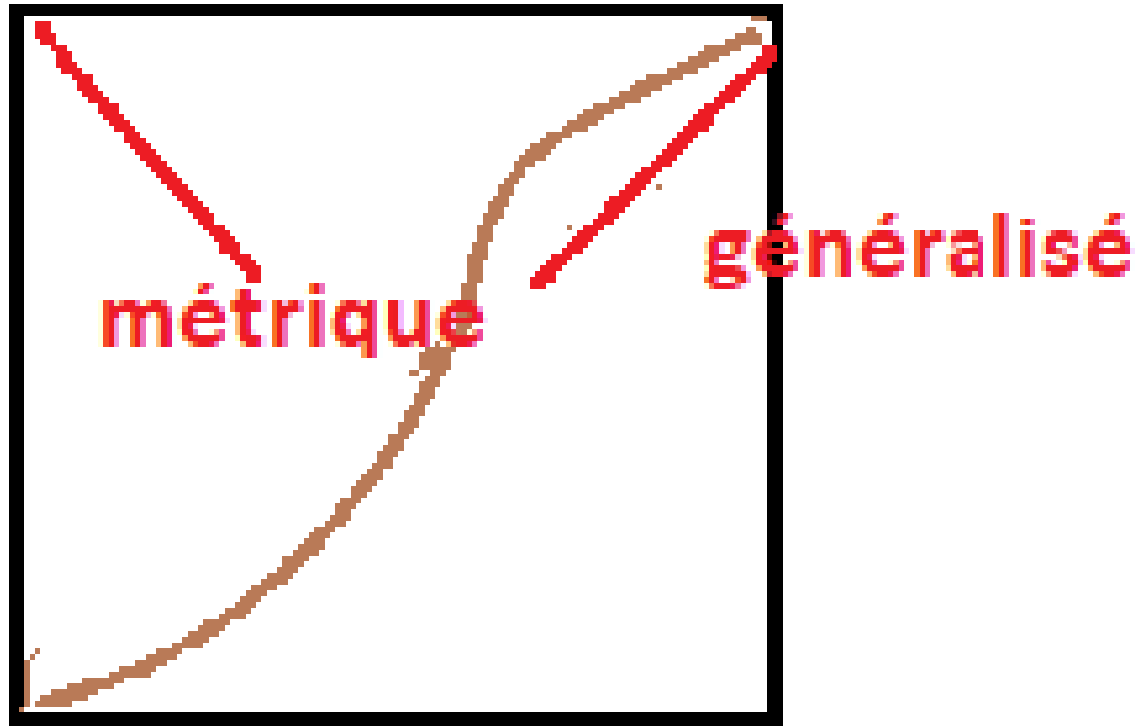


Illustration simpliste de la différence entre calage métrique et calage généralisé

Une condition nécessaire évidente pour que ‘ça’ marche encore pour le calage généralisé :

G (défini uniquement par Z) atteint tous les sommets du contour apparent de C_n (défini uniquement par X).

Je conjecture que cette condition est également suffisante.

Elle est elle-même entraînée par $z_k = q_k x_k$ où les q_k sont des matrices diagonales inversibles. Autrement dit pour chaque k les coordonnées des z_k et des x_k ont le même signe. (même orthant)

En particulier si toutes les variables du problème sont positives le calage généralisé à poids bornés fonctionne comme un calage métrique.

Ouf, c'est fini !

Merci de votre présence, voire de votre attention