

Méthode d'utilisation d'enquête à un niveau géographique où l'échantillon est faible

François JEGER
Insee

I - LA METHODE UTILISEE

La méthode utilisée pour l'estimation du revenu moyen par ménage pour les bassins d'emploi d'Alsace est celle présentée par Robert FAY et Roger HERRIOT dans "Journal of the American Statistical Association" Juin 1979 sous le titre "Estimates of Income for Small Places : "An Application of James STEIN Procedures to Census Data".

L'idée est la suivante :

le revenu par ménage est connu par l'enquête revenus fiscaux avec des erreurs d'échantillonnage fortes dans les zones géographiques fines. D'autres variables, dites auxiliaires, corrélées avec ce revenu, sont mieux connues à partir de sources exhaustives : par exemple, le revenu par foyer fiscal.

Une régression où le revenu par ménage est expliqué par ces variables auxiliaires fournit des "revenus théoriques" pour chacune de ces zones.

Le revenu par ménage est alors estimé par une moyenne de la valeur observée et de la valeur théorique pondérées par l'inverse de leurs variances respectives. Autrement dit, si l'échantillon dans une zone est faible on tiendra peu compte de la valeur observée et beaucoup de la valeur théorique, si l'échantillon est fort ce sera l'inverse.

a - La formulation mathématique : le James Stein et ses dérivés

Le James-Stein - ou estimateur raccourci - est un estimateur biaisé d'un paramètre dont la variance est plus faible que la moyenne d'un échantillon. Le risque quadratique peut être ainsi réduit. Herriot en propose une extension sous la forme suivante :

Soit $Y_i \sim N(\theta_i, D_i)$ le revenu moyen observé dans la zone i , θ_i étant inconnu (à estimer), D_i la variance liée à la taille de l'échantillon de la zone i .

On suppose par ailleurs que $\theta_i \sim N(X_i \beta, A)$ X_i étant une ligne de variables auxiliaires et A la variance inconnue de l'explication de θ_i par les variables auxiliaires.

Si on suppose que les $Y_i - \theta_i$ (erreurs d'échantillonnage) et les $\theta_i - X_i \beta$ (erreurs "d'explication" des θ_i par les X_i) sont indépendants.

$Y \sim N(X \beta, A + D)$ où $A + D$ diagonale de terme $A + D_i$

L'estimation des moindres carrés généralisés fournit un revenu théorique

$$(1) \quad Y_i^* = X_i \hat{\beta} = X_i (X^t (A + D)^{-1} X)^{-1} X (A + D)^{-1} Y$$

L'estimateur retenu est

$$\theta_i^* = \frac{A Y_i + D_i Y_i}{A + D_i}$$

A est alors déterminé par les résidus, puisque

$$E \left(\sum_{i=1}^n \frac{(y_i - y_i^*)^2}{A + D_i} \right) = k - p$$

où k est le nombre d'observations

p est le nombre de paramètres explicatifs

y_i le revenu observé

y_i^* le revenu théorique

Pour des raisons asymptotiques, A sera défini par l'égalité

$$(2) \quad \sum_{i=1}^n \frac{(y_i - y_i^*)^2}{A + D_i} = k - p$$

Il faudrait pour calculer y_i dans (1) connaître la valeur de A donnée par l'équation (2) qui suppose y_i^* déjà connu, pour rompre ce cercle vicieux une méthode itérative est possible.

b - Résolution itérative des équations 1 et 2

La valeur Y_i^* dépend de la variance A , posons $Y_i^* = Y_i^*(A_n)$ où A_n est un estimateur convergent de A .

Chaque régression donne une série de résidus $Y_i - Y_i^*(A_n)$ où A_n est l'estimation de A détenue par la précédente régression.

Ces résidus vont permettre de définir un meilleur estimateur A_{n+1} de A de la façon suivante :

$$\begin{aligned} \text{en posant } f(A_n) &= \frac{(Y_i - Y_i^*(A_n))^2}{A_n + D_i} \\ g(A_n) &= \frac{(Y_i - Y_i^*(A_n))^2}{(A_n + D_i)^2} \end{aligned}$$

$A_{n+1} = A_n + (k - p - f(A_n)) / g(A_n)$ (si cette valeur est positive, sinon $A_{n+1} = 0$)

La suite A_n converge vers une valeur A telle que $k - p = f(A) = 0$ (c'est-à-dire l'équation 2) et on initialise par $A_0 = 0$

Remarque : la convergence est rapide (moins de 10 itérations) elle se justifie par l'approximation par g de la dérivée de f par rapport à A .

$f'(A_n) (A_{n+1} - A_n) + f(A_n) = k - p$ définit un algorithme convergent si f est contractante.

II - APPLICATION A L'EVALUATION DU REVENU MOYEN PAR MENAGE PAR BASSIN D'EMPLOI EN ALSACE

Le revenu moyen par ménage est connu par l'enquête revenus fiscaux de 1979 avec un faible échantillonnage au niveau du bassin d'emploi (une dizaine d'observations dans certains bassins d'emploi).

Soit Y_i ce revenu moyen observé et D_i l'erreur d'échantillonnage associée.

Les variables auxiliaires sont tirées de la banque de données locales.

- 1
 X_i = revenu imposable par foyer fiscal
(source exhaustive)
- 2
 X_i = nombre moyen de foyers fiscaux par ménage
- 3
 X_i = taille moyenne du ménage
- 4
 X_i = % d'agriculteurs dans le bassin d'emploi
- 5
 X_i = % d'étrangers dans le bassin d'emploi
- 6
 X_i = % de la population âgée de + de 60 ans

$Y_i \sim N(\theta_i, D_i)$ θ_i à estimer

$\theta_i = X_i \beta + u_i$ avec $u_i \sim N(0, A)$ A à estimer.

La méthode présente conduit à calculer un revenu théorique

$$Y_i^* = X_i (X^t (A+D)^{-1} X)^{-1} X^t (A+D)^{-1} Y$$

ou $A + D$ est la matrice diagonale de terme général $A + D_i$ dont on calcule la moyenne avec Y_i pondérée par l'inverse des variances

$$\Theta_i^* = \frac{Y_i^* D_i + Y_i A}{A + D_i} \text{ est l'estimateur du revenu moyen par ménage}$$

Il est possible de transformer la régression des moindres carrés généralisés en une régression des moindres carrés ordinaires pondérés. En effet puisque la matrice $A + D$ est diagonale en posant :

$$\sigma_i = \frac{1}{\sqrt{A + D_i}}$$

la régression devient :

$$Y = X \beta + u \text{ avec } u \sim N(0, \sigma^{-2})$$

où σ est la matrice diagonale donc le terme général est σ_i soit encore

$$\sigma Y = \sigma X \beta + v \text{ avec } v \sim N(0, I)$$

L'estimateur de β est donc le même que celui de la régression des moindres carrés ordinaires où les observations sont pondérées par les T_i .

La procédure itérative permettant de déterminer à la fois A et Y_i^* (décrite en b) est donc simplifiée puisque :

$$f(A) = \frac{\sum_i u_i^2}{A + D_i} = \sum_i (\sigma_i u_i)^2 = \sum_i \hat{v}_i^2$$

et

$$g(A) = - \frac{\sum_i u_i^2}{(A + D_i)^2} = - \sum_i (\sigma_i u_i)^2 = - \sum_i \hat{v}_i^2$$

Le passage d'une itération à une suivante est donc définie par

$$A_{n+1} = A_n + (k - p - f(A_n)) / g(A_n).$$

III - LES RESULTATS

a - Régression sur 90 départements français

Le revenu moyen par ménage (REME) est expliqué dans un modèle linéaire en fonction du revenu moyen par foyer fiscal (RF84),

du nombre de foyers par ménage (FOYME), du % d'agriculteurs (PAGRI), le % d'étrangers (PETR), le % de la population de + 60 ans (PAGE), la taille des ménages (TAMEM).

Sur 90 observations (départements) pondérées par le carré du nombre d'observations de l'enquête revenu 79 (pour tenir compte de la variance d'échantillonnage) le R2 est 0,95.

La régression est alors :

$$\begin{aligned} \text{REME} = & - 6.799 + 0,769 \text{ RF84} + 45.407 \text{ FOYME} + 8.509 \text{ PETR} \\ & (- 0,47) \quad (13,45) \quad (4,5) \quad (0,626) \\ & - 11.449 \text{ TAMEM} - 57.510 \text{ PAGRI} - 7.808 \text{ PAGE} + u \\ & (- 2,9) \quad (- 2,25) \quad (- 0,42) \quad (\text{R2} = 0,95) \end{aligned}$$

Les chiffres entre parenthèses sont les students de chaque paramètre.

Les résidus suivent une loi normale avec une probabilité de 0,85, les "students" de PETR et PAGE n'étant pas significativement différents de 0, la proportion d'agriculteurs n'étant pas un indicateur démographique et pouvant avoir un effet en Alsace opposé à celui de la France entière, nous retiendrons 3 variables auxiliaires :

la régression devient alors :

$$\begin{aligned} \text{REME} = & - 24.850 + 0,89 \text{ RF84} + 51.530 \text{ FOYME} - 11.303 \text{ TAMEN} + u \\ & (- 2,9) \quad (16,1) \quad (4,8) \quad (3,7) \quad \text{R2} = 0,94 \end{aligned}$$

La taille du ménage étant corrélée au nombre de foyers par ménage, nous retiendrons seulement deux variables auxiliaires : le revenu par foyer fiscal et le nombre de foyers fiscaux par ménage, la régression est alors :

$$\begin{aligned} \text{REME} = & - 23.037 + 1,054 \text{ RF84} + 16.995 \text{ FOYME} + u \\ & (- 2,5) \quad (28,6) \quad (2,92) \quad \text{R2} = 0,93 \end{aligned}$$

Le R2 diminue peu d'une régression à l'autre, les deux dernières variables auxiliaires sont donc suffisantes à l'explication de REME. De plus, dans la dernière équation le coefficient de RF84 est vraisemblable (si on admet un coefficient d'actualisation 1,35 entre 79 et 84 et que l'abattement des revenus imposables est 0,3).

Nous calculerons sur la base de la dernière régression les estimations de revenus de bassin d'emploi en Alsace. Les estimations résultant des autres régressions sont présentées en annexe.

b - Estimations des revenus des bassins d'emploi

Si THEO est le revenu théorique d'un bassin d'emploi, c'est-à-dire calculé à partir de la dernière régression et VARTH sa

variance associée (mesure de l'erreur de la régression), si RMEEF est le revenu observé dans l'enquête revenu (faible échantillon) et VAREMP sa variance associée (mesure de l'erreur de l'échantillonnage), la méthode d'HERRIOT conduit à retenir l'estimateur du revenu moyen.

$$\text{ESTIM} = \frac{\text{THEO} \times \text{VAREMP} + \text{RMEEF} \times \text{VARTH}}{\text{VAREMP} + \text{VARTH}}$$

(moyenne pondérée des revenus théoriques et observés par l'inverse de leur variance)

dont le coefficient de variation est :

$$\text{CVEST} = \left(\frac{\text{VAREMP} \times \text{VARTH}}{\text{VARTH} + \text{VARTH}} \right)^{1/2}$$

Le tableau suivant donne ces estimations de revenu moyen par ménage par bassin d'emploi en Alsace. On trouvera aussi en annexe les estimations résultant d'un choix de régression différent.

Estimation du revenu moyen par ménage par base d'échantillon New
 10:29 WEDNESDAY, SEPTEMBER

DES	VIL	THEO	VARTH	RNEEF	VAREMP	ESTIM	CVEST	INF	SUP
1	WISSEMBO	2275.3	1422509	70121	83092533	62543.3	1251.50	60070.0	65015.9
2	NIEDERDR	2075.3	1242798	59717	43150053	62309.3	1099.09	59855.0	64163.5
3	HAUEHAU	513.1	1250702	57260	61634374	66557.3	1111.53	64373.7	68736.0
4	SARRE-UN	54292.3	903357	64704	241143599	54331.4	951.55	52466.4	56195.5
5	SAVERNE	60201.1	935283	63615	40463632	60340.9	931.23	53417.7	62264.1
6	STRASBOU	63930.7	125690	78641	39530204	69011.0	353.97	68317.2	69704.8
7	MOLSHHEIM	67917.9	994934	72694	89306538	67970.5	991.98	66026.2	69914.8
8	SCHIRMEC	55343.3	371045	47737	32437865	55309.4	607.77	54116.1	56500.6
9	SELESTAT	61636.1	455004	63824	79340658	61648.5	672.63	60330.2	62966.8
10	STE-MARI	55335.2	23504	47035	107302666	55313.4	531.75	54271.2	56355.7
11	COLMAR	66444.5	395488	76872	50339413	66541.3	626.42	65313.5	67769.1
12	HF-BRISS	23389.4	1113290	75607	97119251	68473.9	1051.46	66413.0	70534.7
13	GUEBAILL	64570.2	200518	73323	54594365	64705.2	888.23	62964.2	66446.3
14	THANN-CE	62739.7	677520	60138	29294290	62680.6	813.94	61085.5	64276.2
15	MULHOUSE	67127.5	157511	72213	13036793	67257.4	394.50	66494.2	68030.6
16	ST-LOUIS	76155.7	295234	55540	51975448	76096.0	541.82	75034.0	77159.0
17	ALT-KIRCH	55335.3	1153296	59055	168794662	66285.9	1070.59	64187.5	68384.3

↑ Coefficient de variabilité de ESTIM

↑ Intervalle de confiance à 95%

ESTIMATION (moyenne pondérée de THEO et RNEEF)

des revenus moyen par ménage en 79

↑ Variance empirique

↑ Revenu observé (faible échantillon)

↑ Variance Théorique

↑ Revenu par ménage théorique expliqué par: - revenu par foyer - nombre de foyers par ménage

↑ Base d'échantillon

c - Régression sur 34 bassins d'emploi. Calcul de la variance théorique par itération

Pour déterminer un revenu théorique fonction du revenu par foyer fiscal et des autres variables auxiliaires, la méthode d'Herriot est appliquée sur 34 bassins d'emploi d'Alsace et Lorraine. La variance théorique est calculée suivant la méthode décrite en I.b. Cinq itérations suffisent pour obtenir la valeur de la variance théorique A.

Nbre d'itérations	R^2 de la régression sur 34 bassins d'emploi REME = f(RF84, PAGE, PAGRI, TAMEN)	A	$\left[\frac{\sum (y_i - v_i)^2}{k-p} \right]$ A + Di
1	0,445	8.10^6	12
2	0,442	13.10^6	4,09
3	0,440	$14,5.10^6$	0,68
4	0,440	$14,57.10^6$	0,04
5	0,440	$14,57.10^6$	0,002

Le R^2 de la régression ne s'améliore pas d'une itération à l'autre, par contre l'égalité

$$\frac{\sum (y_i - v_i)^2}{A + D_i} = 34 - 4$$

est rapidement vérifiée avec une variance théorique de $14,57.10^6$. Cette valeur est supérieure à celle issue de la régression sur 90 départements. Cela provient en partie du faible nombre de bassins d'emplois.

Les régressions sont successivement :

itération 1 : REME = - 61.282 + 1,097 RF84 - 169 PAGE + 302 PAGRI + 20.267 TAMEM

itération 2 : REME = - 58.431 + 1,091 RF84 - 213 PAGE + 300 PAGRI + 20.505 TAMEM

itération 3 : REME = - 57.171 + 1,08 RF84 - 232 PAGE + 298 PAGRI + 20.267 TAMEM

itération 4 : REME = - 56.924 + 1,087 RF84 - 236 PAGE + 298 PAGRI + 20.221 TAMEM

itération 5 : REME = - 56.907 + 1,087 RF84 - 236 PAGE + 298 PAGRI + 20.218 TAMEM

Les estimations de revenus par bassin d'emploi d'Alsace et de Lorraine sont présentées à la page suivante.

En ajoutant à ces 34 observations, 13 bassins d'emploi de Franche-Comté, on n'améliore pas le R^2 de la régression (0,40) et la variance théorique est 17.10 6.

ANNEXE

Estimation des revenus moyens par ménages en 79. Régression sur les 34 bassins d'Alsace et de Lorraine. Variance théorique calculée par situation

OBS	DE	THEO	VARTH	REME	VAREMP	ESTIM	CVEST	INP	SUP
1	4111	53510.2	14579517	53354.2	17936618	53269.1	2836.93	50708.3	61829.6
2	4112	53229.9	14579517	52932.5	40260013	54415.7	3271.02	49001.3	50920.1
3	4113	53747.6	14579517	61590.2	10157233	61336.7	3469.74	57091.0	60532.7
4	4121	54095.7	14579517	42055.7	25295133	52250.4	3041.14	46295.8	53217.1
5	4122	53761.0	14579517	59120.7	15931503	64300.0	2795.91	50314.7	59786.5
6	4125	63904.0	14579517	54731.5	374492347	64410.4	3779.32	56285.1	71832.7
7	4130	53442.4	14579517	63325.5	11752813	65649.1	2551.53	60543.1	70657.1
8	4141	55029.9	14579517	59630.9	29045944	67233.3	3115.63	61120.7	73339.9
9	4152	65107.9	14579517	59325.0	54356122	61932.5	3199.35	55711.6	58253.3
10	4150	53901.5	14579517	61372.9	20403634	59673.1	2915.05	54157.7	55589.6
11	4160	52935.2	14579517	44000.1	27152930	50512.9	3030.15	44475.8	56550.0
12	4171	60255.9	14579517	65200.8	49149463	61024.4	3353.23	55052.1	69196.7
13	4172	52043.9	14579517	52531.2	37732396	52197.1	3243.47	45839.9	58554.3
14	4130	60031.6	14579517	57050.3	37752575	59237.1	1243.11	52980.6	55593.6
15	4191	56543.4	14579517	52745.1	12229006	54523.3	2576.24	42463.7	59577.9
16	4192	62155.6	14579517	52545.3	138036363	64647.5	3631.36	56930.1	71165.0
17	4193	54725.7	14579517	65135.9	24950024	57643.7	3034.19	51820.5	63797.7
18	4271	67153.9	14579517	79121.9	93359761	67554.3	3551.13	60594.6	74515.0
19	4272	65309.0	14579517	59717.6	42337712	64262.2	3293.11	57797.9	70726.5
20	4273	70292.2	14579517	57260.6	67004106	73320.2	3460.37	66543.9	80103.5
21	4274	61053.4	14579517	64704.4	223597793	61236.3	3699.61	54035.1	68537.5
22	4275	52591.1	14579517	63615.3	33798604	62564.2	3254.33	56435.7	69242.7
23	4276	67095.6	14579517	78641.4	43637672	69997.1	3305.21	63507.7	76466.4
24	4277	70071.3	14579517	72694.9	83631346	70972.1	3523.53	64060.0	77973.2
25	4278	53270.7	14579517	47737.3	74155699	52369.7	3490.59	45528.3	59211.4
26	4279	63373.9	14579517	53524.9	70371257	63450.9	3477.36	56035.3	70266.5
27	4281	50497.0	14579517	47095.4	106022800	50080.0	3531.31	44067.2	57105.0
28	4282	64305.6	14579517	73872.3	43174219	68073.7	3345.49	61510.0	74530.9
29	4283	79217.4	14579517	75507.7	95040433	72763.9	3555.35	71735.4	85732.4
30	4284	65005.0	14579517	73338.0	52005342	67731.1	3374.57	61117.1	74745.1
31	4285	64535.6	14579517	60135.6	25216153	52924.7	3039.44	56967.4	69882.0
32	4286	68814.6	14579517	72213.2	13039072	70605.4	2626.23	65453.0	75752.9
33	4237	70267.7	14579517	65540.5	50661572	73870.5	3354.74	67275.6	80465.4
34	4289	70336.0	14579517	59056.0	149584829	69334.2	3644.53	62190.3	75479.0

Lorraine

Alsace

ANNEXE

Estimation basée sur la régression : 90 départements, 6 variables (RF84, TAMEM, PAGRI, PETR, PAGE et FOYME)

OBS	NCC	THEO	VARTH	RMEEF	VAREMP	ESTIM	CVEST	INF	SUP
1	WISSEMO	64346.5	2634779	70121	93096533	64524.0	1528.03	61391.3	67556.2
2	NIEDERBR	5757.9	1933719	59717	43150053	65509.6	1362.11	62839.9	68179.4
3	HAGUENAU	6934.2	1633364	67266	61684374	70395.3	1231.02	67894.5	72906.1
4	SARRE-UN	5371.1	3111510	64704	241142599	58909.1	2501.35	53418.4	64399.7
5	SAVERNE	51219.3	2235267	63615	40453532	61347.9	1470.97	58464.9	64231.0
6	STRASBOU	72534.1	950069	75641	39580204	72677.2	953.22	70799.3	74565.1
7	MOLSHEIM	63359.3	1433368	72694	89306658	68437.4	1139.79	66105.4	70769.3
8	SCHIRMEC	63437.5	2435225	47737	92437265	62937.9	1539.01	59971.4	66004.3
9	SELESTAT	51764.9	1355603	63324	79340658	61799.0	1154.54	59536.1	64061.9
10	STE-MARI	63455.4	2390890	47025	107302666	63100.1	1526.17	60108.3	66091.4
11	COLMAR	6794.7	1770055	73872	50333419	69812.9	1307.66	67249.9	72375.9
12	NF-BRISS	63355.2	3672817	75807	97119251	63808.9	1881.22	60121.7	67496.1
13	SUESWILL	67950.1	2327083	73835	54594365	68192.9	1493.97	65264.7	71121.1
14	THANN-CE	63426.3	1537331	60133	29294290	67985.0	1245.26	65547.3	70428.7
15	MULHOUSE	70322.0	529532	72213	13036793	70395.8	713.35	68997.7	71794.0
16	ST-LOUIS	75360.7	742069	65540	51975448	75705.5	855.36	74029.0	77382.0
17	ALTIRICH	66013.3	3553774	59055	163794662	65677.7	1813.40	62323.5	69432.0