

# La coordination d'échantillons d'enquêtes auprès des entreprises

---

Fabien Guggemos, Olivier Sautory  
*INSEE*



*Journées de Méthodologie Statistique 2012 - Paris*



# Numéros aléatoires

---

Technique des numéros aléatoires utilisées par de nombreux INS pour coordonner des échantillons auprès des enquêtes entreprises.

A chaque unité  $k$  de la population  $U$  (y compris aux nouvelles unités) est attribué, indépendamment des autres, un nombre  $\omega_k$ , tiré aléatoirement selon une loi uniforme sur  $[0,1 [$ .

**Tirage de Poisson**: sélection des unités  $k$  dont les numéros aléatoires sont dans l'intervalle  $[d, d + \pi_k[$ ,  $\pi_k =$  probabilité d'inclusion de l'unité  $k$ .

**Tirage aléatoire simple (TAS)** : sélection des  $n$  unités ayant les plus petits numéros aléatoires  $\omega_k$  supérieurs à  $d$ .

## Définition d'une fonction de coordination

---

Fonction de coordination  $g$  = fonction borélienne (i.e. mesurable) de  $[0,1[$  vers  $[0,1[$ , déterministe, qui préserve la probabilité uniforme :

Si  $P$  est la probabilité uniforme sur  $[0,1[$ , alors la probabilité image doit vérifier:

$$P^g = P.$$

→ Notamment, pour tout intervalle  $I = [a,b [$  inclus dans  $[0,1[$ :

$$P(g^{-1}(I)) = P^g(I) = P(I) = b - a$$

→ Et en particulier...

De même que les  $(\omega_k)_{k \in U}$ , les numéros aléatoires transformés  $(g(\omega_k))_{k \in U}$  sont indépendants et suivent une loi uniforme sur  $[0,1[$

# Sélection des unités

---

Pour *chaque* unité  $k$  dans le champ d'intérêt :

- un numéro aléatoire  $\omega_k$  dit permanent, (invariant dans le temps)
- une fonction de coordination  $g_{k,t}$  qui évolue à chaque nouveau tirage d'échantillon  $t = 1, 2, \dots$

## 1. Tirage de Poisson

Sélection des unités  $k$  telles que  $g_k(\omega_k) \in [0, \pi_k [$   
où  $\pi_k$  = probabilité d'inclusion de l'unité  $k$

$$P(k \in S) = P(g_k(\omega_k) \in [0, \pi_k [) = P^{g_k}([0, \pi_k [) = P([0, \pi_k [) = \pi_k$$

Et les tirages sont indépendants.

# Sélection des unités

---

## 2. Tirage aléatoire simple stratifié (TASST)

A l'intérieur d'une strate, sélection des  $n$  unités  $k$  ayant les  $n$  plus petits numéros aléatoires  $g_k(\omega_k)$ .

Et comme les nombres transformés  $g_k(\omega_k)$  sont indépendants et suivent la loi uniforme  $P...$

→ Les  $n$  plus petits nombres  $g_k(\omega_k)$  fournissent donc bien un échantillon aléatoire simple de taille  $n$  dans la strate considérée.

$$p(S) = \binom{N}{n}^{-1} \quad (N : \text{taille de la population})$$

# Une procédure “pas à pas” pour prendre en compte les charges de réponse

---

$\omega = (\omega_1, \dots, \omega_k, \dots, \omega_N)$  = vecteur des numéros aléatoires permanents attribués aux unités de la population U.

$I_{k,t}(\omega)$  = **Indicatrice**, égale à 1 si l'unité k est sélectionnée dans l'échantillon t, et 0 sinon :  $k \in S_t \Leftrightarrow I_{k,t}(\omega) = 1$

$\gamma_{k,t}$  = **charge de réponse** pour l'entreprise k lors de l'enquête t.  
La **charge réelle de réponse** est aléatoire :  $\gamma_{k,t}(\omega) = \gamma_{k,t} I_{k,t}(\omega)$

**Charge cumulée** pour l'unité k :  $\Gamma_{k,t}(\omega) = \sum_{u \leq t} \gamma_{k,u} \cdot I_{k,u}(\omega)$

Principe : définir des fonctions  $g_{k,t}$  de sorte que :  
 $\Gamma_{k,t-1}(\omega_{(1)}) < \Gamma_{k,t-1}(\omega_{(2)}) \Rightarrow g_{k,t}(\omega_{k,(1)}) < g_{k,t}(\omega_{k,(2)})$

## Deux difficultés

---

1. Substituer à  $\Gamma_{k,t}(\omega)$  une fonction de  $\omega_k$  seulement, notée  $\Gamma'_{k,t}(\omega_k)$ , qui fournit une bonne approximation de  $\Gamma_{k,t}(\omega)$ .

Tirage de Poisson :  $I_{k,t}(\omega)$  dépend seulement de  $\omega_k$  ( $I_{k,t}(\omega) = I_{k,t}(\omega_k)$   
 $= 1$  si  $g_{k,t}(\omega_k) \in [0, \pi_k[$ , 0 sinon)  
 $\rightarrow \Gamma_{k,t}(\omega)$  dépend seulement de  $\omega_k$ , noté  $\Gamma'_{k,t}(\omega_k)$ .

TASST:  $I_{k,t}(\omega)$  dépend “essentiellement” de  $\omega_k$ , mais pas seulement... On approxime  $I_{k,t}(\omega)$  par  $I'_{k,t}(\omega) = E(I_{k,t}(\omega) | \omega_k)$ , puis  $\Gamma_{k,t}(\omega)$  par  $\Gamma'_{k,t}(\omega_k)$ .

2. Définir la fonction de coordination  $g_{k,t}$  de sorte que :

$$\Gamma'_{k,t-1}(\omega_{k,1}) < \Gamma'_{k,t-1}(\omega_{k,2}) \Rightarrow g_{k,t}(\omega_{k,1}) < g_{k,t}(\omega_{k,2})$$

# Construction d'une fonction de coordination

---

Pour un critère général de coût  $C(\omega_k)$ , comme la charge de réponse cumulée approximée : plus il est petit, plus on souhaite que la probabilité de sélection de l'unité  $k$  dans l'échantillon  $t$  soit élevée.

$$C(\omega_1) < C(\omega_2) \Rightarrow g_C(\omega_1) < g_C(\omega_2) \quad (1)$$

Si  $C$  borélienne, on peut construire une première fonction  $G_C$  :

$$G_C(\omega) = P^C (]-\infty, C(\omega)[) = P(C^{-1}]-\infty, C(\omega)[) = P(u|C(u) < C(\omega))$$

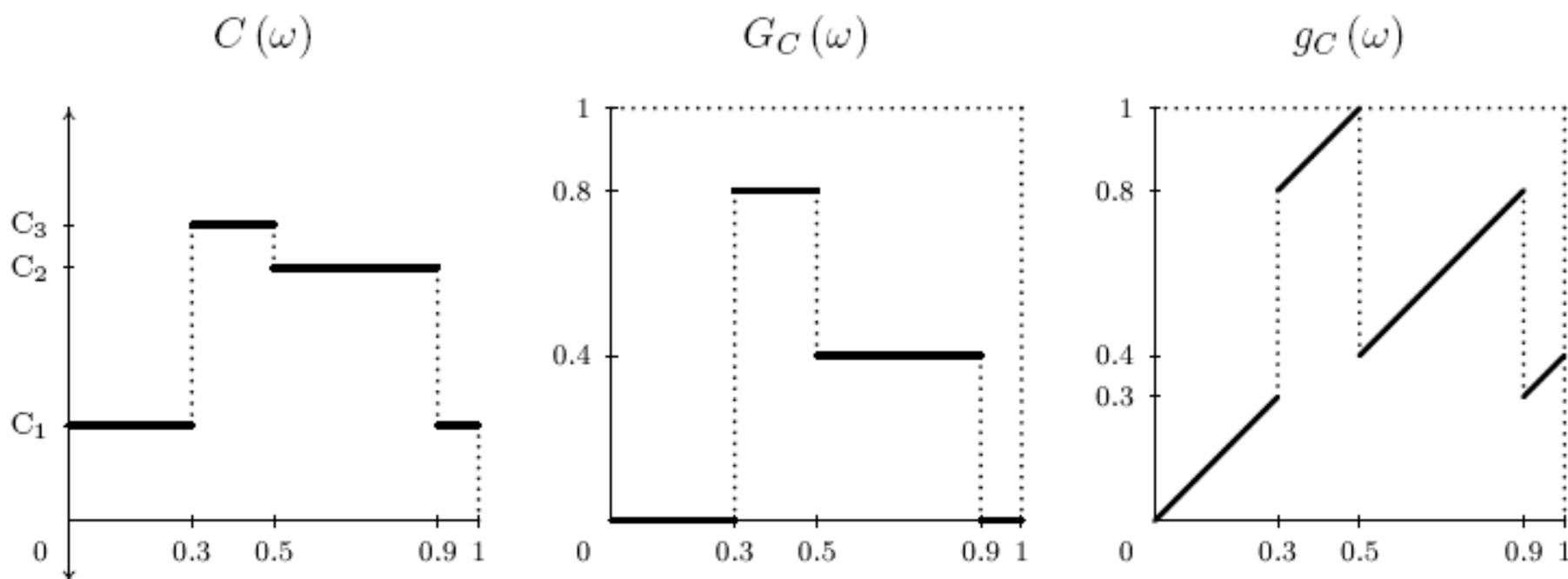
$$"G_C = F_C(C)"$$

Si  $C$  admet des paliers (boréliens de  $[0;1[$  de mesure de Lebesgue non nulle sur lesquels  $C$  est constante), alors  $G_C$  admet des paliers et ce sont les mêmes que pour  $C$ .

Si pas de palier,  $G_C$  est une fonction de coordination vérifiant (1) !!!

## Construction d'une fonction de coordination

Et si  $C$ , donc  $G_C$ , ont des paliers ? Une simple transformation affine des paliers permet de modifier  $G_C$  en une fonction de coordination  $g_C$  satisfaisant (1) !!!



On peut ensuite introduire des critères de coûts secondaires...

# Application au Tirage de Poisson

---

Initialisation :  $\Gamma_{k,0}(\omega_k) = 0$     $g_{k,1}(\omega_k) = \omega_k$

$$I_{k,1}(\omega_k) = \mathbb{I}_{[0, \pi_{k,1}[}(\omega_k) \quad \Gamma_{k,1}(\omega_k) = \gamma_{k,1} \mathbb{I}_{[0, \pi_{k,1}[}(\omega_k)$$

Pour l'échantillon  $S_t$ , on choisit une fonction de coordination  $g_{k,t}$  associée à l'unité  $k$ .

$$k \in S_t \Leftrightarrow g_{k,t}(\omega_k) \in [0, \pi_{k,t}[$$

Si l'on définit :  $A_{k,t} = g_{k,t}^{-1}[0, \pi_{k,t}[$

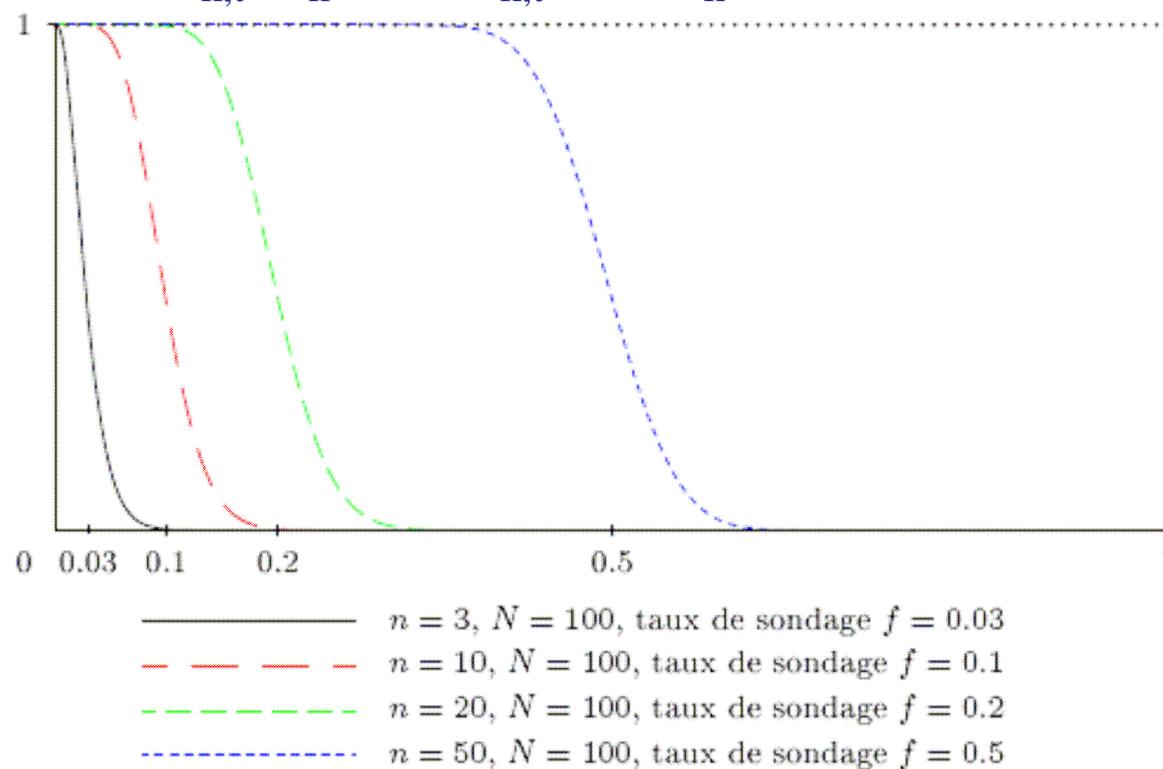
→ indicatrice de sélection :  $I_{k,t}(\omega_k) = \mathbb{I}_{A_{k,t}}(\omega_k)$

**Pour la sélection d'un échantillon « séparé »  $S_t$  :**

critère  $C_{k,t}(\omega_k) =$  charge cumulée  $\Gamma_{k,t-1}(\omega_k)$ , on en déduit  $g_{k,t}$

# Application au TASST

## 1. Approximation $\Gamma'_{k,t}(\omega_k) = E(I_{k,t}(\omega) | \omega_k)$ ,



2. Implémentation informatique :  $\Gamma'_{k,t}$  et charges cumulées espérées *approximées* par des fonctions constantes par morceaux ( $[0;1[$  subdivisé en  $L$  sous-intervalles). → Puis procédure similaire au Tirage de Poisson

# Première simulation, cas simple

---

Population fixe dans le temps de taille 100.

Charge initiale = 0 pour chaque unité.

Sélection de 20 échantillons (SAS).

Pour chaque échantillon : taille  $n = 25$ , charge de réponse = 1,

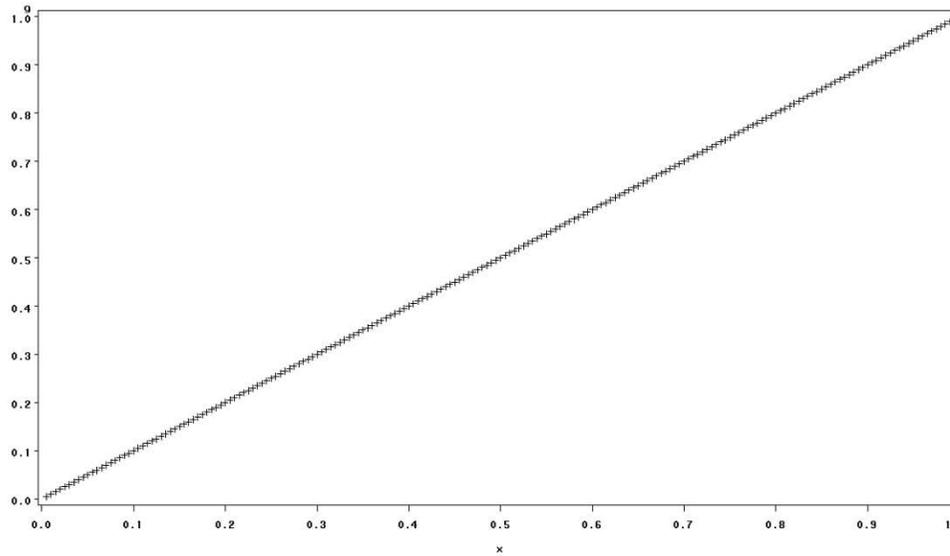
*sauf :*

Echantillons n°3 et n°15 : taille  $n = 50$ , charge de réponse = 3

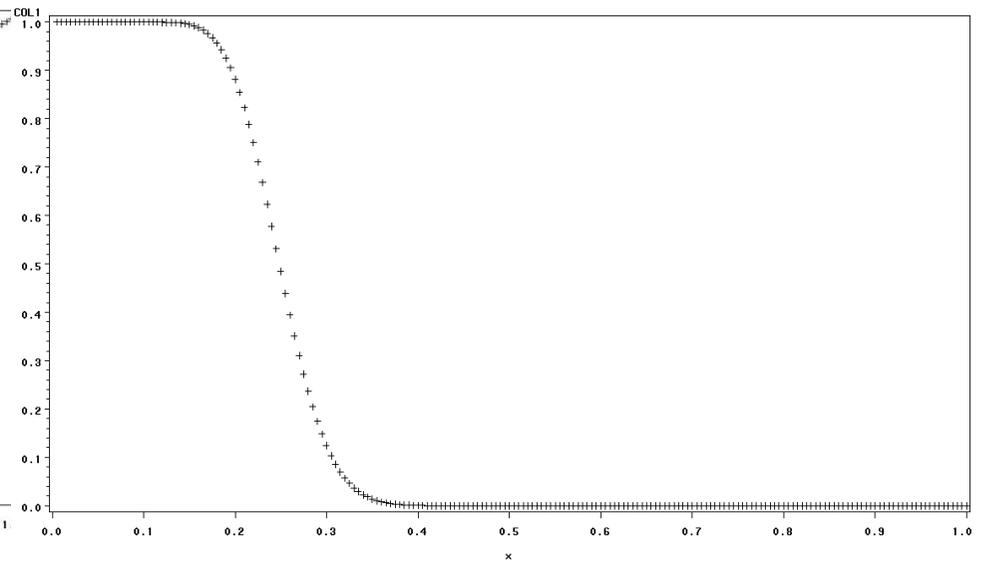
Echantillons n°10 et n°11 : taille  $n = 10$ , charge de réponse = 2

→ Charge cumulée moyenne espérée = 7.4

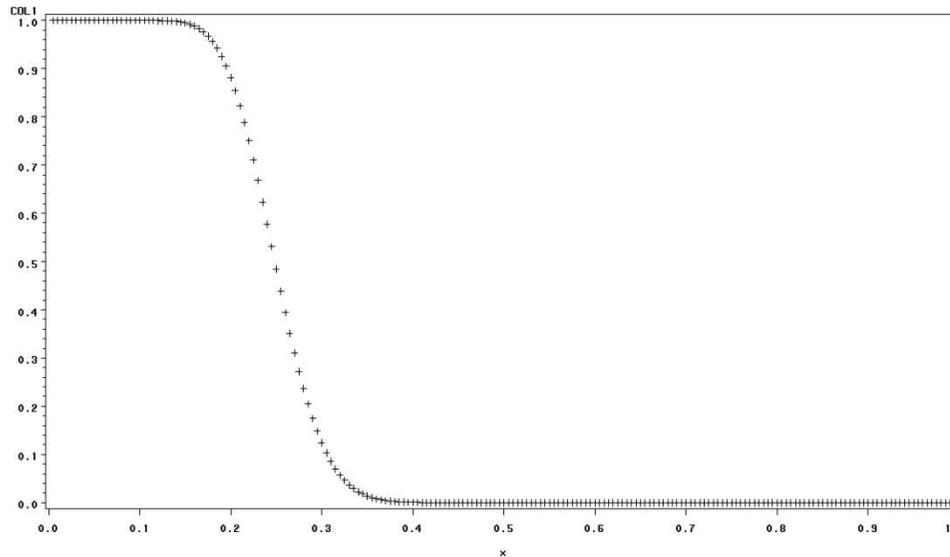
Survey 1, samplesize = 25, response burden = 1  
Coordination function



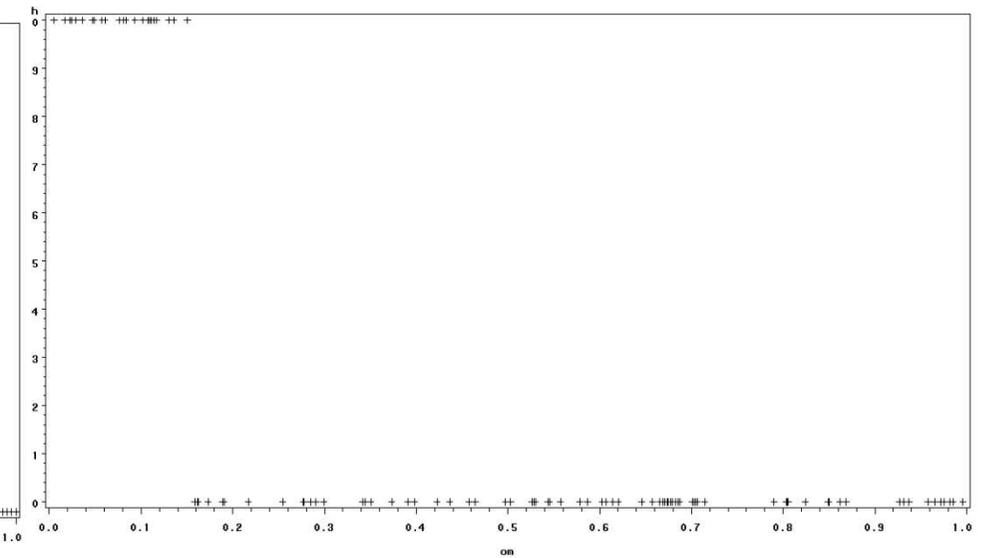
Survey 1, samplesize = 25, response burden = 1  
Approximate indicator function



Survey 1, samplesize = 25, response burden = 1  
Cumulative burden after sampling

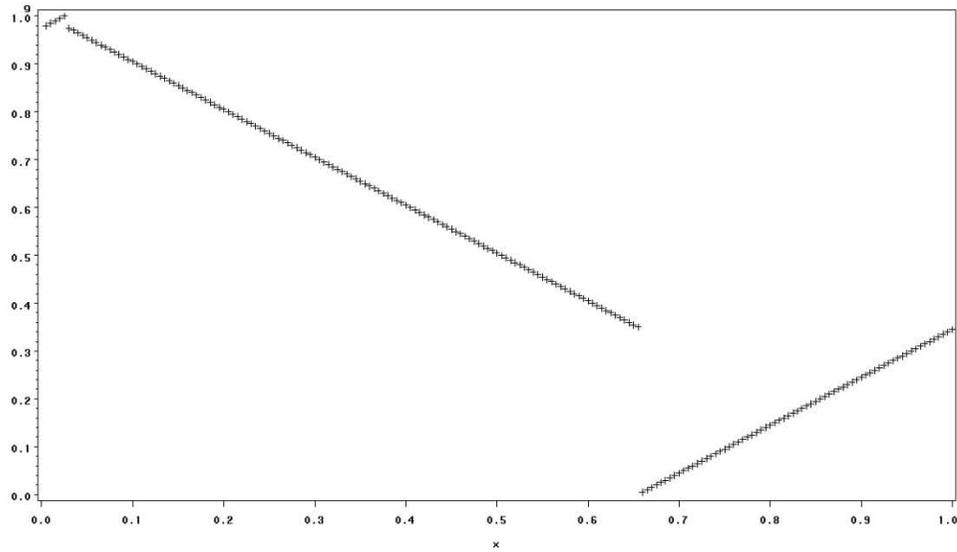


Survey 1, samplesize = 25, response burden = 1  
Sample



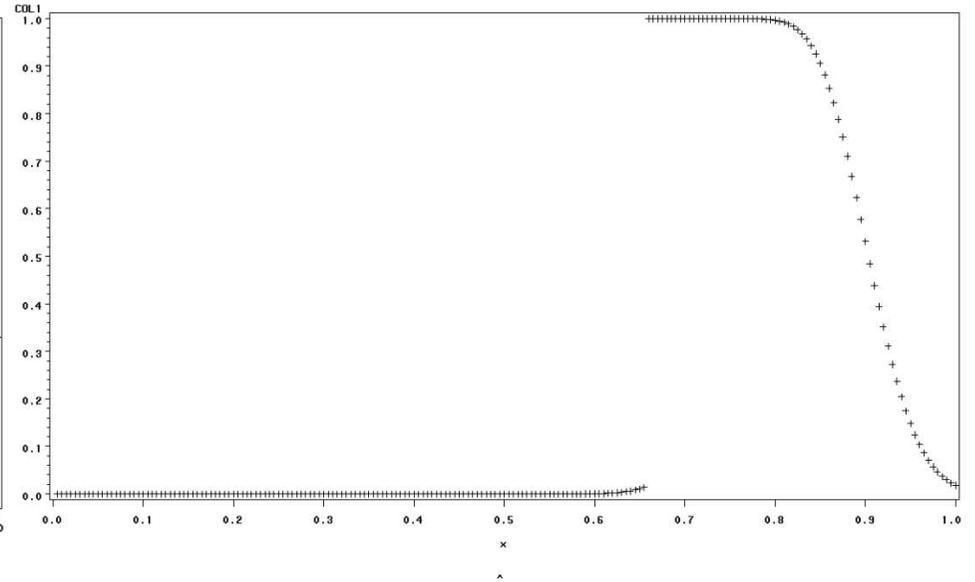
Survey 2, samplesize = 25, response burden = 1

Coordination function



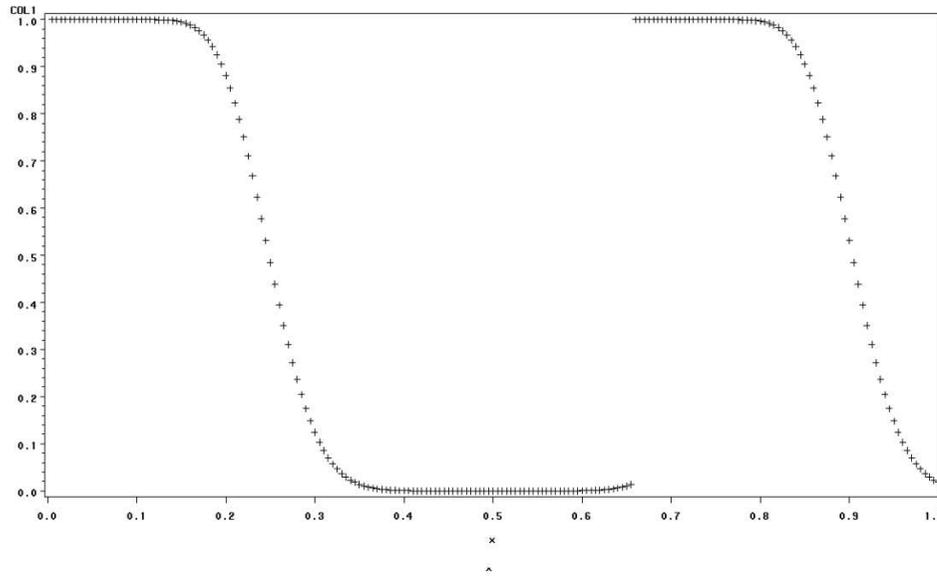
Survey 2, samplesize = 25, response burden = 1

Approximate indicator function



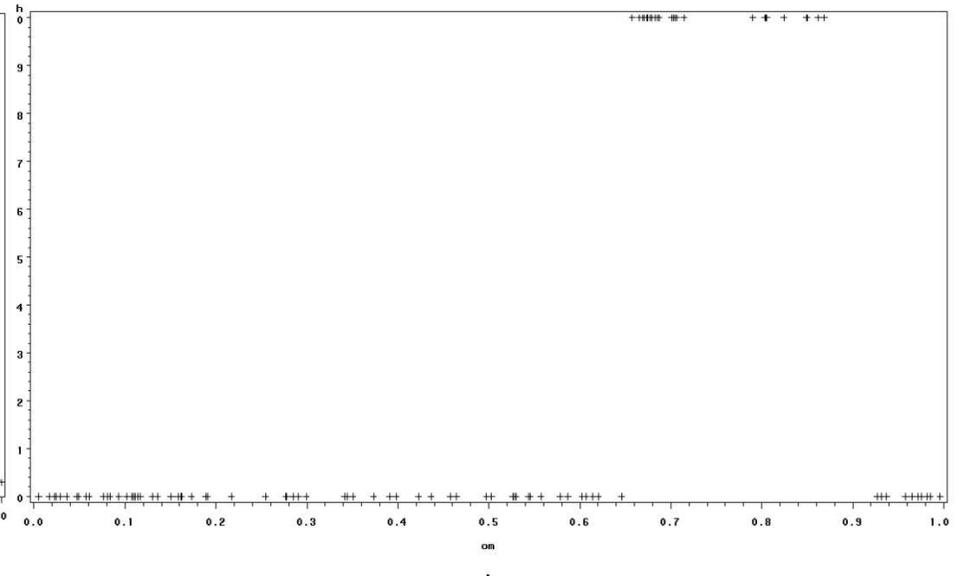
Survey 2, samplesize = 25, response burden = 1

Cumulative burden after sampling

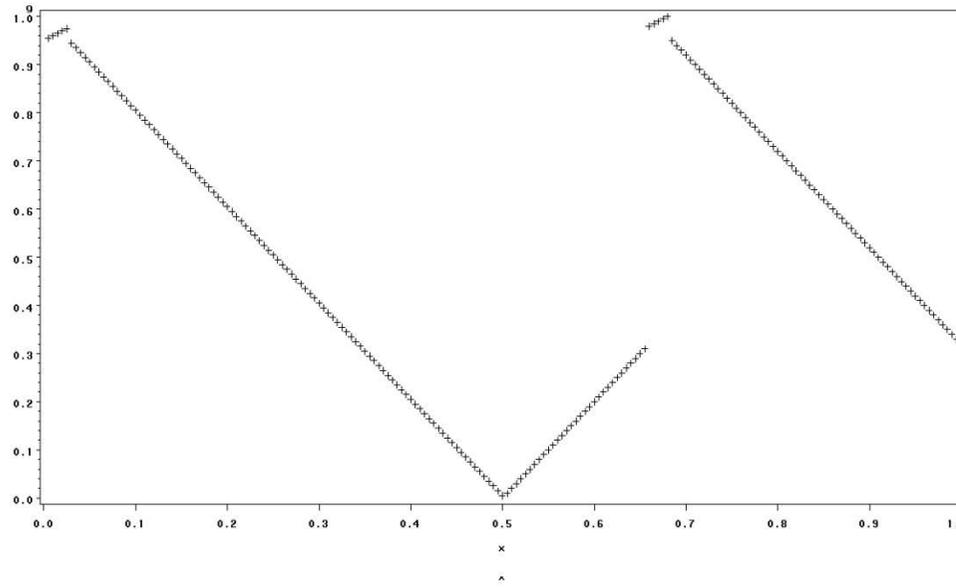


Survey 2, samplesize = 25, response burden = 1

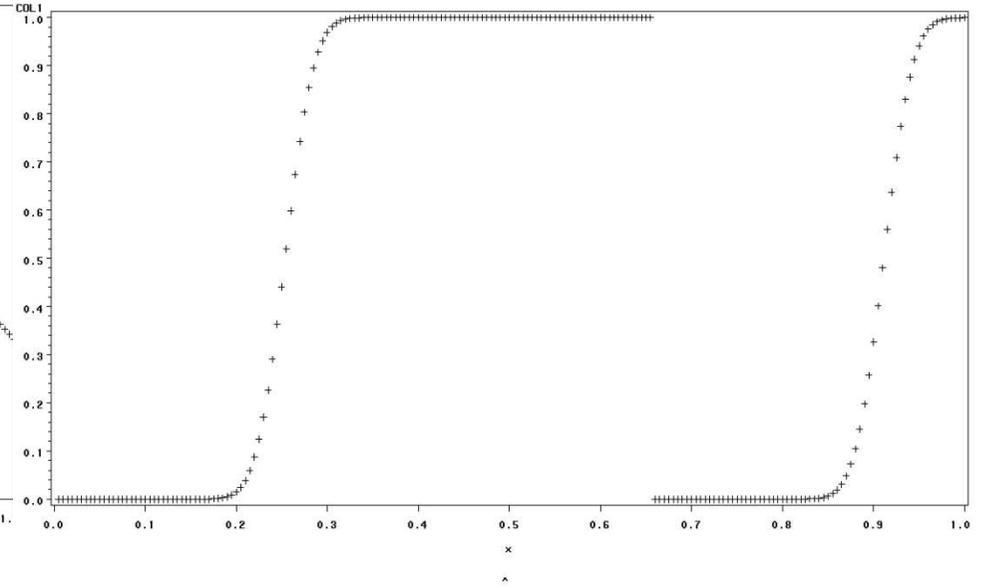
Sample



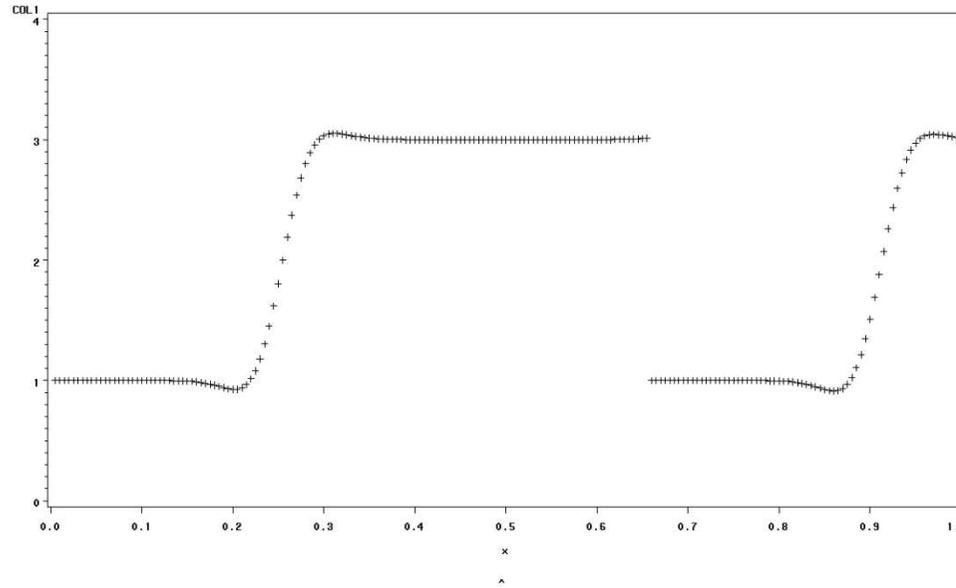
Survey 3, samplesize = 50, response burden = 3  
Coordination function



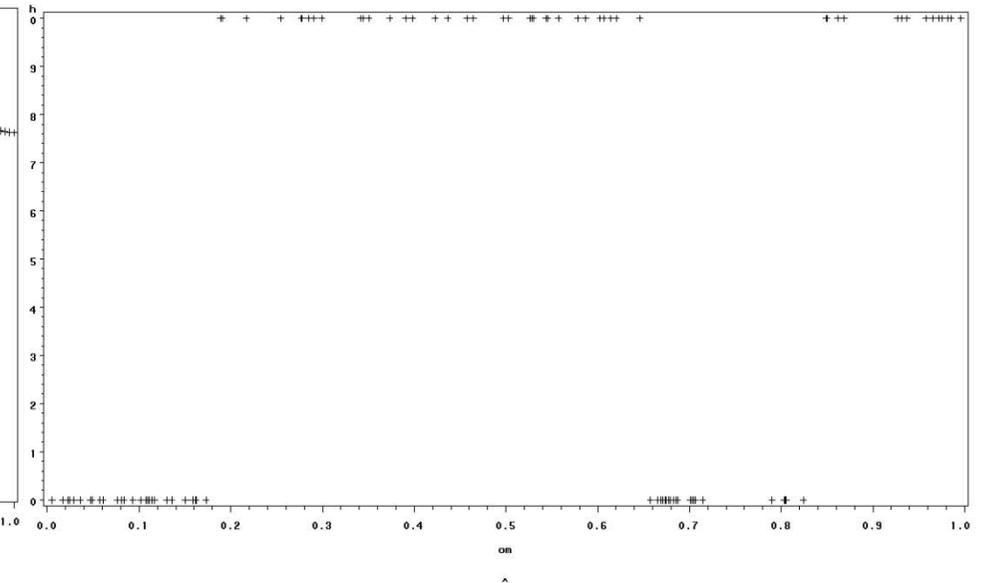
Survey 3, samplesize = 50, response burden = 3  
Approximate indicator function



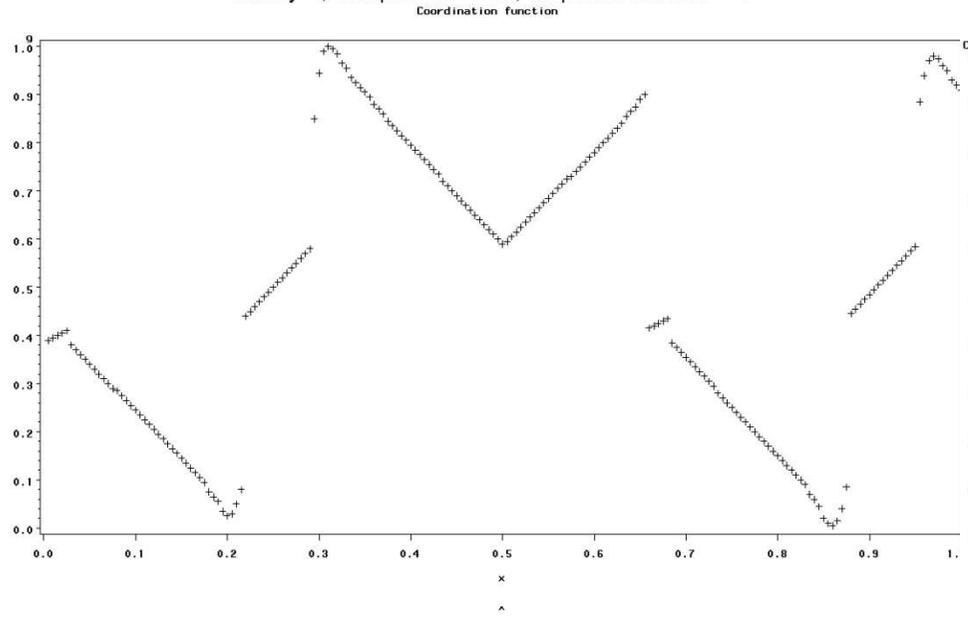
Survey 3, samplesize = 50, response burden = 3  
Cumulative burden after sampling



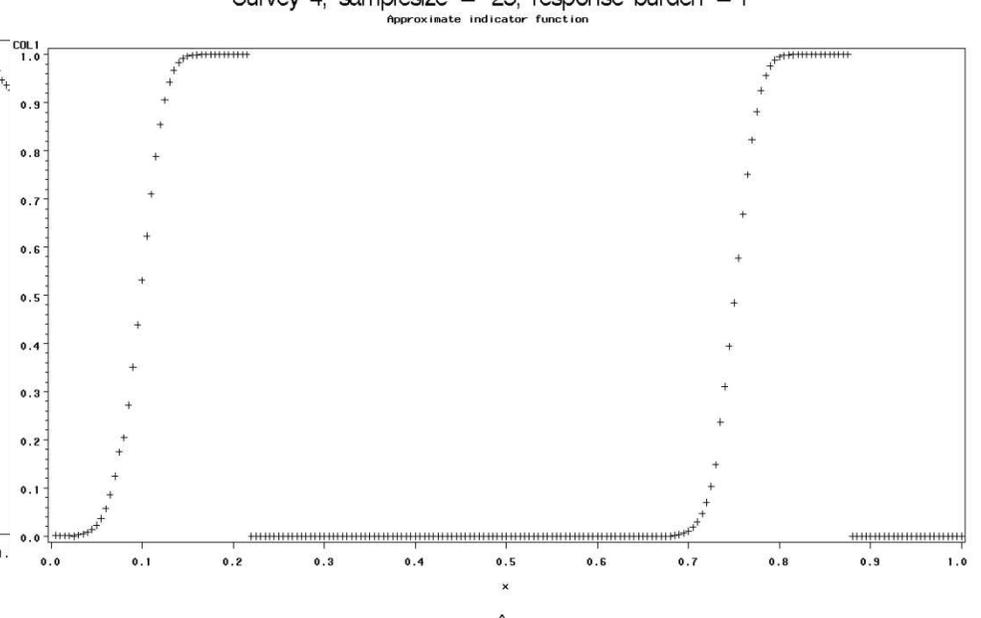
Survey 3, samplesize = 50, response burden = 3  
Sample



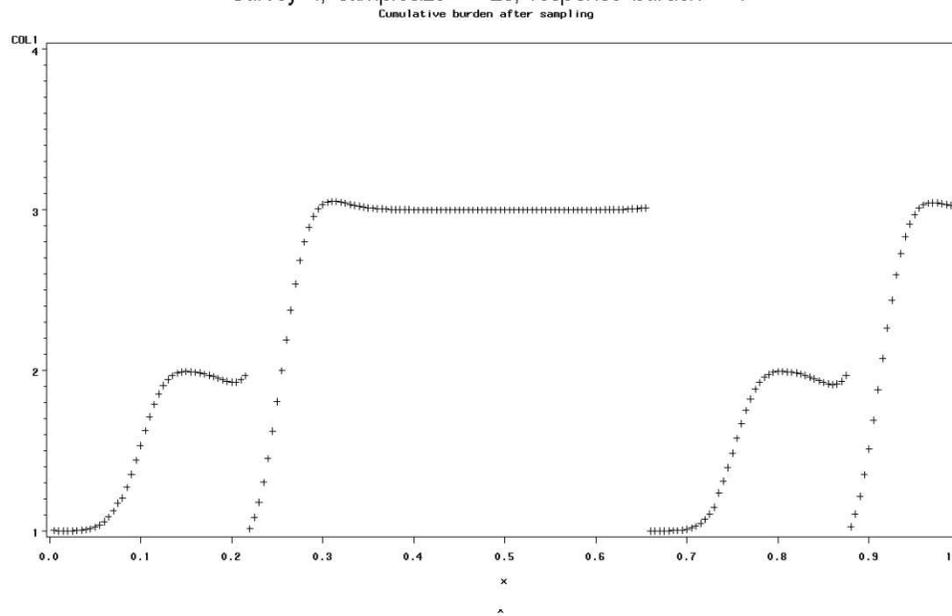
Survey 4, samplesize = 25, response burden = 1



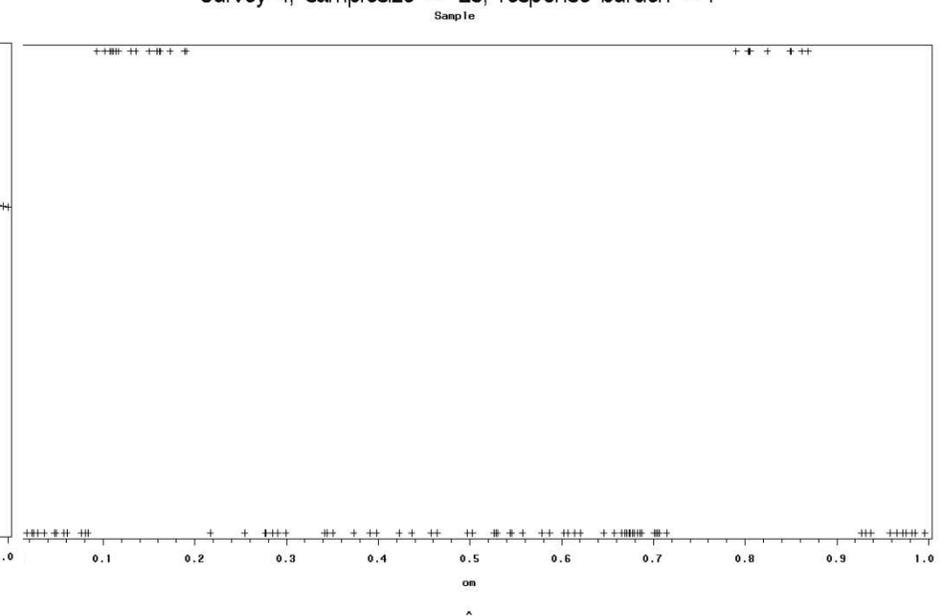
Survey 4, samplesize = 25, response burden = 1



Survey 4, samplesize = 25, response burden = 1

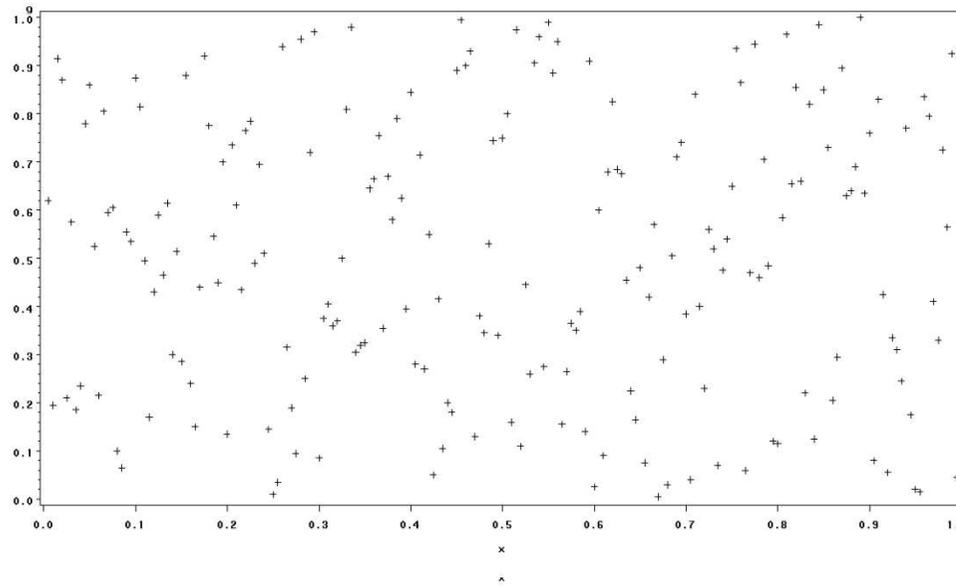


Survey 4, samplesize = 25, response burden = 1



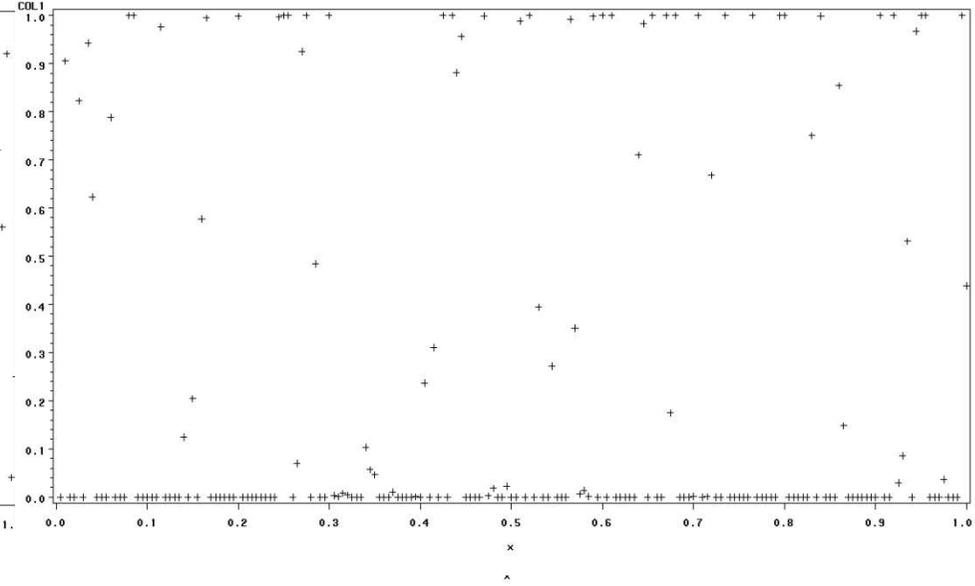
Survey 20, samplesize = 25, response burden = 1

Coordination function



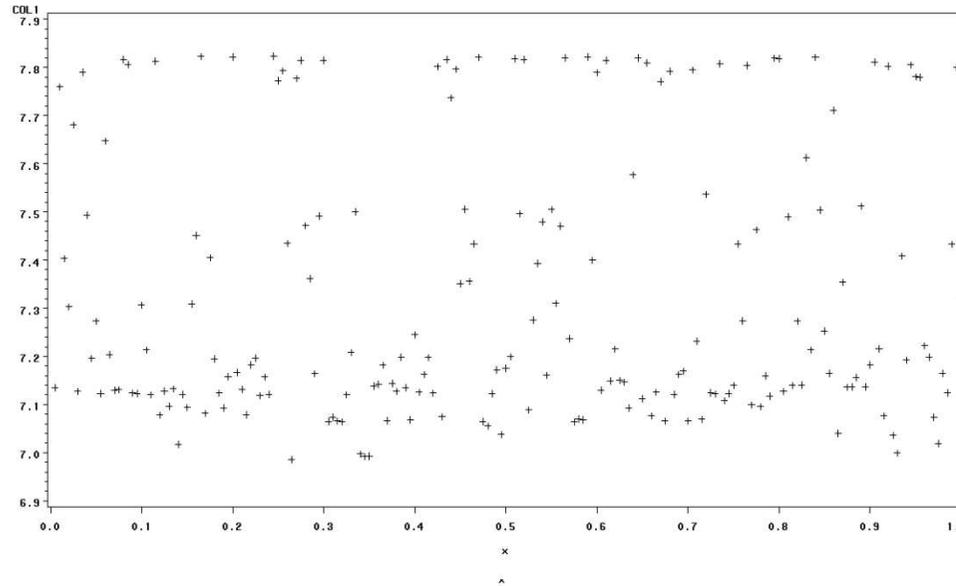
Survey 20, samplesize = 25, response burden = 1

Approximate indicator function



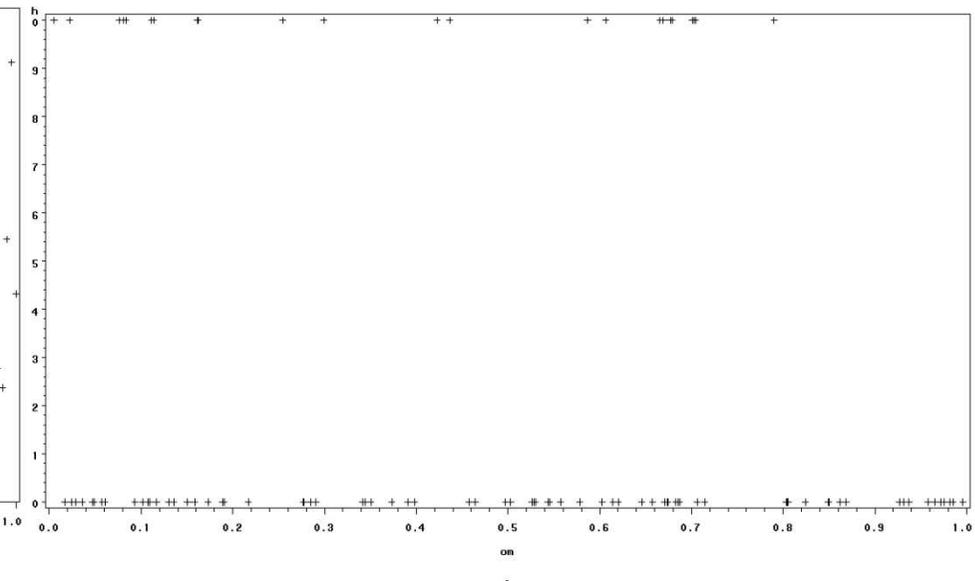
Survey 20, samplesize = 25, response burden = 1

Cumulative burden after sampling



Survey 20, samplesize = 25, response burden = 1

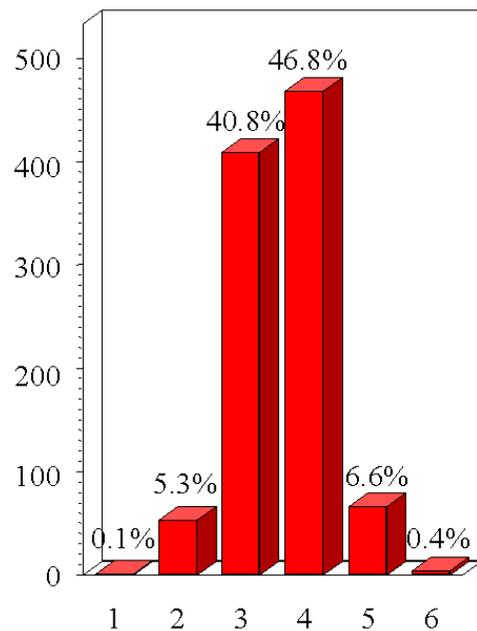
Sample



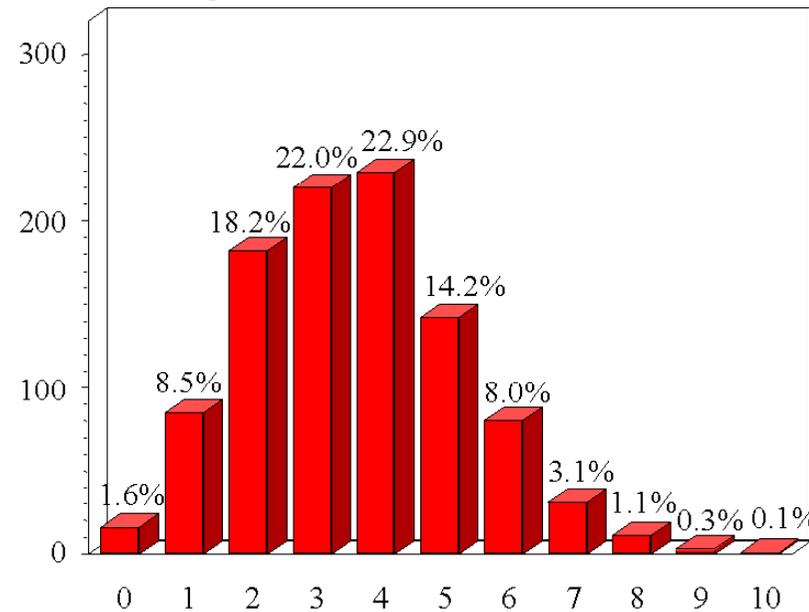
# Simulations : performances de la coordination en regard de divers paramètres

Tirage de 20 échantillons successifs (TASST) sur une population évoluant dans le temps.

Nombre d'unités en fct de la fréquence d'échantillonnage



Cas coordonné



Cas indépendant

2 critères de performance :

$$R_{\sigma} = \frac{\sigma_{\text{Coord}}}{\sigma_{\text{Indep}}} < 1$$

$$\Delta_s = S_{\text{Coord}} - S_{\text{Indep}} < 0$$

## Simulations : performances de la coordination en regard de divers paramètres

---

- Globalement, **résultats très positifs** : pour quasi toutes les simulations envisagées,  $R_{\sigma} < 0.5$  (Plus de 50% de gains en termes de répartition de la charge) et  $\Delta_S < 0$
- **Robustesse de la méthode** vis-à-vis de situations instables : population mouvante, stratif. et taux de sondage très variables d'un tirage à l'autre...
- Importance de coordonner négativement une enquête avec suffisamment d'enquêtes du passé : avec au moins  $N_{\text{coord}}$  enquêtes, où  $N_{\text{coord}} > \bar{f}^{-1} - 1$  ( $\bar{f}$  : taux moyen de sondage)

Beaucoup plus de résultats dans l'article associé à cette présentation ...

---

# Merci de votre attention !

## Contact

M. Fabien Guggemos

Tél. : 00 33 1 41 17 50 18

Mail : [fabien.guggemos@insee.fr](mailto:fabien.guggemos@insee.fr)

M. Olivier Sautory

Tél. : 00 33 1 41 17 50 82

Mail : [olivier.sautory@insee.fr](mailto:olivier.sautory@insee.fr)

## Insee

18 bd Adolphe-Pinard  
75675 Paris Cedex 14

[www.insee.fr](http://www.insee.fr)  

Informations statistiques :

[www.insee.fr](http://www.insee.fr) / Contacter l'Insee

09 72 72 4000

(coût d'un appel local)

du lundi au vendredi de 9h00 à 17h00