**IVS**

**Institut für
Volkswirtschaftslehre
und Statistik**

No. 586-00

**Survey sampling: a linear game**

Horst Stenger, Siegfried Gabler,
Jochen Schmidt

# Beiträge zur angewandten Wirtschaftsforschung

**Universität Mannheim
A5, 6
D-68131 Mannheim**

# Survey sampling: a linear game

Horst Stenger, Siegfried Gabler, Jochen Schmidt

**Summary**: A linear game consists of two subsets of a vector space with a scalar product. The idea is that players 1 and 2 select, independently, elements of the first and second set, respectively. Then, player 2 has to pay to player 1 the value of the scalar product of the selected elements.

We will discuss survey sampling within the framework of linear games with the statistician in the role of player 2. The vector space to be considered is the set of all symmetric matrices of order $N \times N$ with a scalar product identical with the usual mean squared error. The subset from which the statistician's selection is to be made is neither convex nor compact. Standard results of the theory of linear games have to be modified appropriately.

The existence of minimax strategies will be established. At the same time we hope to improve our understanding of random selection and of the duality between the fixed population approach and model based approaches to the theory of survey sampling.

# 1  Estimating vectors, designs and strategies

Let
$$y_1, y_2, \ldots y_N \in \mathbb{R}$$
be unknown values of a characteristic of interest for units $1, 2, \ldots N$. Suppose a size $n$ sample
$$s = \left\{ i_1, i_2, \ldots i_n \right\} \subset \left\{ 1, 2, \ldots N \right\}$$
is selected and a linear function
$$\sum_{i \in s} a_{si}\, y_i$$
of the observed values $y_i$, $i \in s$ is used to estimate the total $y = \sum y_i$. Define
$$a_{si} = 0 \quad \text{for} \quad i \notin s$$
and consider the *estimating vector* $\underline{a}_s = (a_{s1}, a_{s2}, \ldots a_{sN})'$. With $\underline{y} = (y_1, y_2, \ldots y_N)'$ we have
$$\sum_{i \in s} a_{si}\, y_i = \underline{a}_s{}'\underline{y}$$
and, $\underline{y}$ given, the loss resulting from using $\underline{a}_s$ is

$$\begin{aligned}
\left( \sum a_{si}\, y_i - y \right)^2 &= \left( \underline{a}_s{}'\underline{y} - y \right)^2 \\
&= \left[ \left( \underline{a}_s - \underline{1} \right)'\underline{y} \right]^2 \\
&= \underline{y}'\left( \underline{a}_s - \underline{1} \right)\left( \underline{a}_s - \underline{1} \right)'\underline{y} \\
&= tr\left[ \underline{y}\,\underline{y}'\left( \underline{a}_s - \underline{1} \right)\left( \underline{a}_s - \underline{1} \right)' \right].
\end{aligned}$$

Note that

$$\underline{y}\,\underline{y}'$$
$$\left( \underline{a}_s - \underline{1} \right)\left( \underline{a}_s - \underline{1} \right)'$$

are elements of the vector space of all symmetric $N \times N$-matrices and by the trace operator $tr$ a scalar product is defined in this space.

A *design p* is a probability distribution on the set of all size $n$ samples. An *estimator $\underline{a}$* associates an estimating vector $\underline{a}_s$ with each size $n$ sample $s$. The performance of the *strategy $(p, \underline{a})$, $p$* a design and $\underline{a}$ an estimator, is characterized be the *mean squared error (MSE)*

$$MSE(\underline{y}; p, \underline{a}) = \sum_s p_s \left( \sum_i a_{si} y_i - y \right)^2$$

$$= tr\left[ \underline{y}\,\underline{y}' \sum_s p_s \left( \underline{a}_s - \underline{1} \right)\left( \underline{a}_s - \underline{1} \right)' \right]$$

which is the scalar product of $\underline{y}\,\underline{y}'$ and the *risk generating matrix*

$$\sum_s p_s \left( \underline{a}_s - \underline{1} \right)\left( \underline{a}_s - \underline{1} \right)'$$

introduced by CHENG and LI(1983).

Define

$$\Theta_0 = \left\{ \underline{y}\,\underline{y}' : \underline{y} \in \mathbb{R}^N \right\}$$

and let $\mathcal{A}_0$ be the set of all risk generating matrices. Subsequently, the linear game

$$\left( \Theta_0, \mathcal{A}_0 \right)$$

(with scalar product as pay-off) will be modified in a way to describe realistically the statistician's decision problem.

## 2  Superpopulation models

The selection of a strategy will be based on prior information. Suppose there is some evidence of a linear model $M(\Omega)$ behind the vector $\underline{y}$, i.e. $\underline{y}$ may be interpreted as a realization of a vector $\underline{Y}$ of random variables $Y_1, Y_2 \ldots Y_N$ satisfying

$$\underline{Y} = X\,\underline{\beta} + \underline{\varepsilon}$$

with a $N \times K$ matrix $X$ of full rank and a random vector $\underline{\varepsilon} = (\varepsilon_1, \varepsilon_2 \ldots \varepsilon_N)'$ with

$$E\,\underline{\varepsilon} = \underline{0}, \ var\,\underline{\varepsilon} = \Omega.$$

3

On the base of a size $n$ sample $s$ a linear predictor of $\sum Y_i$ may be written as

$$\underline{a}_s{}' \underline{Y}$$

where $\underline{a}_s$ is an estimating vector. Unbiasedness of this predictor is equivalent to

$$E\left(\underline{a}_s{}' \underline{Y} - \underline{1}' \underline{Y}\right) = 0$$

i.e.

$$X'\left(\underline{a}_s - \underline{1}\right) = \underline{0}.$$

Under this condition

$$\begin{aligned} var\left(\underline{a}_s{}' \underline{Y} - \underline{1}' \underline{Y}\right) &= (\underline{a}_s - \underline{1})'\Omega(\underline{a}_s - \underline{1}) \\ &= tr\left(\Omega(\underline{a}_s - \underline{1})(\underline{a}_s - \underline{1})'\right) \end{aligned}$$

is minimized by *best linear unbiased (BLU) regression predictors*. However, usually the model $M(\Omega)$, called *superpopulation model*, is not reliable in a strict sense and, therefore, the use of standard regression predictors is not justified.

The strategy $(p, \underline{a})$ to be used should perform well in „small neighbourhoods"of

$$L = \left\{ X \underline{\beta} : \underline{\beta} \in \mathbb{R}^K \right\}.$$

Let $\underline{\varepsilon}$, in addition to the earlier assumptions, be normally distributed with $\Omega$ regular. Then, the logarithm of the density function of $\underline{Y}$ is a linear function of

$$\left(\underline{y} - X \underline{\beta}\right)'\Omega^{-1}\left(\underline{y} - X \underline{\beta}\right)$$

and a natural distance between $\underline{y}$ and $L$ is defined by

$$\min_{\underline{\beta} \in \mathbb{R}^K} \left(\underline{y} - X \underline{\beta}\right)'\Omega^{-1}\left(\underline{y} - X \underline{\beta}\right) = \underline{y}'\left[\Omega^{-1} - \Omega^{-1}X(X'\Omega^{-1}X)^{-1}X'\Omega^{-1}\right]\underline{y}$$

$$= \underline{y}' U \underline{y}, \text{ say,}$$

with $U$ non-negative definite and of rank $N - K$ satisfying $U X = 0$. Hence, a strategy $(p, \underline{a})$ with risk generating matrix

$$W = \sum p_s\left(\underline{a}_s - \underline{1}\right)\left(\underline{a}_s - \underline{1}\right)'$$

will be excluded from consideration unless

$$tr\, \underline{y}\, \underline{y}'\, W$$

4

is bounded on the *parameter set*

$$\Theta_1 = \left\{ \underline{y}\,\underline{y}' \ : \ \underline{y}'\,U\,\underline{y} \leq 1 \right\}$$

which is the case if and only if $p_s > 0$ implies

$$X'(\underline{a}_s - \underline{1}) = \underline{0}$$

i.e. $(p, \underline{a})$ is *representative* in the sense of HAJEK (1959). Note that the representativity of $(p, \underline{a})$ is equivalent with unbiasedness of $\underline{a}'_s \underline{Y}$ in the model $M(\Omega)$.

# 3  Minimax strategies

Let $\mathcal{A}_x$ be the set of all risk generating matrices of representative strategies. Then, the games

$$(\Theta_1, \mathcal{A}_x)\,, (\widetilde{\Theta}_1, \mathcal{A}_x)$$

are of interest, where $\widetilde{\Theta}_1$ is the convex hull of $\Theta_1$.

**Theorem:** *For the linear game $\left(\widetilde{\Theta}_1, \mathcal{A}_x\right)$ a value $v$ and matrices $\overset{*}{V} \in \widetilde{\Theta}_1$, $\overset{*}{W} \in \mathcal{A}_x$ exist with*

$$tr \ V\overset{*}{W} \leq v \leq tr \ \overset{*}{V}W$$

*for all $V \in \widetilde{\Theta}_1, W \in \mathcal{A}_x$.*

The proof of this theorem is given in the appendix.

Here, we consider a consequence. Obviously, a representative strategy $(\overset{*}{p}, \overset{*}{\underline{a}})$ exists with

$$\overset{*}{W} = \sum \overset{*}{p}_s \ (\overset{*}{\underline{a}}_s - \underline{1})(\overset{*}{\underline{a}}_s - \underline{1})'.$$

Define

$$\Theta = \left\{ \underline{y} : \underline{y}\,\underline{y}' \in \Theta_1 \right\} = \left\{ \underline{y} : \underline{y}'U\underline{y} \leq 1 \right\}.$$

Then, for $\underline{y} \in \Theta$ and $(p, \underline{a})$ representative

$$MSE\left(\underline{y}; \overset{*}{p}, \overset{*}{\underline{a}}\right) \leq v \leq \sup_{\underline{y} \in \Theta} MSE\left(\underline{y}; p, \underline{a}\right)$$

i.e. $(\overset{*}{p},\overset{*}{\underline{a}})$ is minimax on $\Theta$ in the usual sense of survey sampling. If a statistician applies $(\overset{*}{p},\overset{*}{\underline{a}})$ his risk $MSE(y;\overset{*}{p},\overset{*}{\underline{a}})$ is a function of $\underline{y} \in \Theta$ and $v$ is an upper bound; in addition, $v$ is the lowest upper bound he can achieve by any strategy.

Further, for all $s \in S$ with $\overset{*}{p}_s > 0$ and $\underline{a}_s$ representative

$$tr \; \overset{*}{V} (\underline{a}_s - \underline{1})(\underline{a}_s - \underline{1})' \;\; \geq \;\; tr \; \overset{*}{V}\overset{*}{W}$$
$$\geq \;\; tr \; \overset{*}{V} (\overset{*}{\underline{a}}_s - \underline{1})(\overset{*}{\underline{a}}_s - \underline{1})'.$$

Hence, for all $s$ with $\overset{*}{p}_s > 0$

$$\overset{*}{\underline{a}}_s{}' \; \underline{Y}$$

is a BLU predictor of $\sum Y_i$ in the model $M(\overset{*}{V})$

$$\underline{Y} = X \underline{\beta} + \underline{\varepsilon}$$

with $E \underline{\varepsilon} = \underline{0}$, $var \underline{\varepsilon} = \overset{*}{V}$. So, a statistician starting with a superpopulation model $M(\Omega)$, not strictly reliable, has strong motivation to base his decision on a model $M(\overset{*}{V})$ derived from $M(\Omega)$ in the way outlined above.

# 4 Related work

Minimax strategies have been derived for the parameter set

$$\Theta^{(1)} = \left\{ \underline{y} \in \mathbb{R}^N : \sum (y_i - \overline{y})^2 \leq 1 \right\}$$

e.g. by BICKEL and LEHMANN (1981), GABLER (1990). STENGER and GABLER (1996) discuss, more generally,

$$\Theta^{(2)} = \left\{ \underline{y} \in \mathbb{R}^N : \sum \sum d_{ij} (y_i - \overline{y})(y_j - \overline{y}) \leq 1 \right\}$$

with $(d_{ij})$ a positive definite $N \times N$ matrix which must be close to the identity matrix I.

$$\Theta^{(3)} = \left\{ \underline{y} \in \mathbb{R}^N : \sum \left( y_i - \frac{y}{x} x_i \right)^2 \leq 1 \right\}$$

6

is an example of a parameter set depending on a vector $\underline{x} = (x_1, x_2, \ldots x_N)'$ with $x_i > 0\,;\, i = 1, 2, \ldots N$. For this set it is shown by STENGER (1990) that the ratio strategy is minimax in an asymptotic sense. Minimax strategies in the strict sense (of the present paper) are derived by GABLER and STENGER (2000) for the set

$$\Theta = \left\{ \underline{y} \,:\, \sum\sum d_{ij}\left( y_i - \frac{y}{x}\,x_i \right)\left( y_j - \frac{y}{x}\,x_j \right) \leq 1 \right\}$$

under the assumption

$$\underline{x} \quad \text{close to} \quad \underline{1},$$
$$(d_{ij}) \quad \text{close to the identity matrix I.}$$

Note that $\Theta$ may be written as

$$\left\{ \underline{y} \,:\, \underline{y}'\,U\,\underline{y} \leq 1 \right\}$$

with $U\,\underline{x} = \underline{0}$. In the present paper $\Theta$ is supposed to satisfy

$$U\,X = 0$$

for a $N \times K$ matrix $X$ of full rank $K \geq 1$, and no further assumptions are needed.

# 5   Appendix: Proof of the Theorem

Consider the orthogonal projector along $X$:

$$I - X(X'X)^{-1}X' = P, \text{ say.}$$

Obviously, $P' = P$ and
$$PWP = W$$
for all $W \in \Theta_2$. Hence,

$$tr\left(\underline{y}\,\underline{y}'W\right) = tr\left(\underline{y}\,\underline{y}'PWP\right)$$
$$= tr\left(P\,\underline{y}\,(P\,\underline{y})'W\right)$$

7

and
$$\Theta_x := \left\{ \underline{y}\,\underline{y}' \in \Theta_1 \ : \ X'\underline{y} = \underline{0} \right\}$$
is a complete subset of $\Theta_1$. Let
$$\widetilde{\Theta}_x$$
be the convex hull of $\Theta_x$. To prove the Theorem of section 3 we will show that an equilibrium point $\left( \overset{*}{V}, \overset{*}{W} \right)$ exists for the game $\left( \widetilde{\Theta}_x, \mathcal{A}_x \right)$.

Let $\mathcal{Q}$ be the set of all $N \times N$-matrices $Q$ which are symmetric and satisfy the condition
$$X'Q = 0.$$
$\mathcal{Q}$ is a vector space with the scalar product $\langle V, W \rangle = tr\, VW$. With
$$\mathcal{Q}_0 := \left\{ Q \in \mathcal{Q} \ : \ Q \ \text{non-negative definite} \right\}$$
we have, obviously,
$$\widetilde{\Theta}_x, \widetilde{\mathcal{A}}_x \subset \mathcal{Q}_0 \subset \mathcal{Q}$$
and
$$\widetilde{\Theta}_x = \left\{ Q \in \mathcal{Q}_0 \ : \ \langle U, Q \rangle \leq 1 \right\},$$
where $\widetilde{\mathcal{A}}_x$ is the convex hull of $\mathcal{A}_x$.

Note that the set $\widetilde{\mathcal{A}}_x$ is introduced for technical reasons, while the set $\mathcal{A}_x$ which is not convex is of primary interest.

By $\mathcal{Q}_0$ an order relation is defined on $\mathcal{Q}$: For
$$A, B \in \mathcal{Q} \ \text{with} \ A - B \in \mathcal{Q}_0$$
we write
$$A \succeq B.$$

This order relation differs from the relation usually considered in statistical decision theory which is defined by the cone of all vectors (matrices, in our case) with non-negative components.

**Lemma 1:** *Choose $Q \in \mathcal{Q}_0$ with $U \succeq Q$. Then*
$$N\langle U, U \rangle \geq \langle Q, Q \rangle \geq 0.$$

Proof: Let $\lambda_1 \geq \cdots \geq \lambda_N$ and $\mu_1 \geq \cdots \geq \mu_N$ be the eigenvalues of $Q$ and $U$, respectively. Then
$$\langle Q, Q \rangle = \sum_{i=1}^{N} \lambda_i^2 \leq N\lambda_1^2 \leq N\mu_1^2 \leq N \sum_{i=1}^{N} \mu_i^2 = N\langle U, U \rangle. \qquad \diamond$$

8

**Lemma 2:** *Let $(W^k)$ be a sequence in $\mathcal{A}_x$ with $\lim W^k = W$. Then, there exists $\overset{*}{W} \in \mathcal{A}_x$ with*

$$W \succeq \overset{*}{W}.$$

Proof: We have

$$W^k = \sum_{s \in S} p_s^k (\underline{a}_s^k - \underline{1})(\underline{a}_s^k - \underline{1})' \quad \forall\, k \in \mathbb{N}.$$

Obviously, there exists a subsequence of $(p^k)$ such that

$$p_s := \lim_{k \in \mathbb{N}'} p_s^k \geq 0 \text{ and } \sum p_s = 1$$

where $\mathbb{N}'$ is a subset of $\mathbb{N}$. We define

$$S_+ := \{ s \in S : p_s > 0 \}$$

and

$$\overset{*}{W}{}^k := \sum_{s \in S_+} p_s^k (\underline{a}_s^k - \underline{1})(\underline{a}_s^k - \underline{1})' \quad \forall\, k \in \mathbb{N}.$$

Then

$$W^k - \overset{*}{W}{}^k \in \mathcal{Q}_0 \quad \forall\, k \in \mathbb{N}.$$

Again, we choose a subset $\mathbb{N}'' \subseteq \mathbb{N}'$ for which

$$\underline{a}_s := \lim_{k \in \mathbb{N}''} \underline{a}_s^k \quad \forall\, s \in S_+$$

exists. Then

$$\overset{*}{W} := \sum_{s \in S_+} p_s (\underline{a}_s - \underline{1})(\underline{a}_s - \underline{1})' \in \mathcal{A}_x$$

and, by the closedness of $\mathcal{Q}_0$,

$$W - \overset{*}{W} = \lim_{k \in \mathbb{N}''} (W^k - \overset{*}{W}{}^k) \in \mathcal{Q}_0. \hspace{2cm} \diamond$$

**Lemma 3:** *Let $\tilde{W} \in \widetilde{\mathcal{A}}_x$. Then, there exists $W \in \mathcal{A}_x$ with*

$$\tilde{W} \succeq W.$$

9

Proof: We have

$$\tilde{W} = \sum_k \lambda_k \sum_s p_s^k (\underline{a}_s^k - \underline{1})(\underline{a}_s^k - \underline{1})'$$

with $\lambda_k \geq 0$, $\sum_k \lambda_k = 1$. Define

$$p_s := \sum_k \lambda_k p_s^k$$

and

$$\underline{a}_s := \sum_k \frac{\lambda_k p_s^k}{p_s} \underline{a}_s^k.$$

Then

$$W := \sum_s p_s (\underline{a}_s - \underline{1})(\underline{a}_s - \underline{1})' \in \mathcal{A}_x.$$

Since the mapping

$$\underline{y} \mapsto \underline{y}\,\underline{y}' \quad \forall \underline{y} \in \mathbb{R}^N$$

is convex with respect to $\succeq$,

$$\tilde{W} \succeq W. \qquad\qquad \diamond$$

Note that, for $V \in \mathcal{Q}_0$, the mapping

$$Q \mapsto \langle V, Q \rangle \quad \forall Q \in \mathcal{Q}_0$$

is increasing with respect to $\succeq$. As a consequence of Lemma 3, only the elements of $\mathcal{A}_x$ are of interest, because every $\tilde{W} \in \widetilde{\mathcal{A}}_x$ is dominated (in the sense of the game) by a risk generating matrix $W \in \mathcal{A}_x$.

With

$$\mathcal{U}_\alpha := \big\{ Q : \alpha\, U \succeq Q \big\}$$

for $\alpha \in \mathbb{R}$, the set

$$\big\{ \alpha : \mathcal{U}_\alpha \cap \widetilde{\mathcal{A}}_x \neq \emptyset \big\}$$

is non-empty. As a consequence of $\mathcal{U}_0 \cap \widetilde{\mathcal{A}}_x = \emptyset$

$$\alpha_0 := \inf \big\{ \alpha : \mathcal{U}_\alpha \cap \widetilde{\mathcal{A}}_x \neq \emptyset \big\}$$

is non-negative. There exists a hyperplane, defined by $Q \in \mathcal{Q} \backslash \{0\}$ and $k \in \mathbb{R}$, separating the interior of $\mathcal{U}_{\alpha_0}$ (relative to $\mathcal{Q}$) and $\widetilde{\mathcal{A}}_x$, i.e.

$$\langle Q, R \rangle \leq k \quad \text{for all} \quad R \in \mathcal{U}_{\alpha_0}$$
$$\langle Q, \tilde{W} \rangle \geq k \quad \text{for all} \quad \tilde{W} \in \widetilde{\mathcal{A}}_x$$
$$k = \sup_{R \in \mathcal{U}_{\alpha_0}} \langle Q, R \rangle.$$

Suppose $Q \notin \mathcal{Q}_0$. Then $\underline{z} \in \mathbb{R}^N$ exists with $X' \underline{z} = \underline{0}$ and $\underline{z}' Q \underline{z} < 0$, i.e.

$$\langle Q, \underline{z}\,\underline{z}' \rangle < 0.$$

Now, for all $\tau > 0$,
$$\tau \,\underline{z}\,\underline{z}' \in \mathcal{Q}_0$$

and therefore
$$-\tau \,\underline{z}\,\underline{z}' \in \mathcal{U}_{\alpha_0}.$$

However,
$$\langle Q, -\tau \,\underline{z}\,\underline{z}' \rangle = -\tau \langle Q, \underline{z}\,\underline{z}' \rangle$$

goes to $\infty$ for $\tau \to \infty$ because of $\langle Q, \underline{z}\,\underline{z}' \rangle < 0$. This is in contradiction to the separating property of the hyperplane considered. Hence,

$$Q \in \mathcal{Q}_0 \backslash \{0\}.$$

$U$ is an element of the interior of $\mathcal{Q}_0$ (relative to $\mathcal{Q}$) and

$$\langle Q, U \rangle > 0.$$

Dividing
$$k = \sup_{R \in \mathcal{U}_{\alpha_0}} \langle Q, R \rangle = \alpha_0 \langle Q, U \rangle$$

by $\langle Q, U \rangle$ we obtain

$$\alpha_0 = \frac{k}{\langle Q, U \rangle} = \sup_{R \in \mathcal{U}_{\alpha_0}} \frac{\langle Q, R \rangle}{\langle Q, U \rangle}$$
$$= \sup_{R \in \mathcal{U}_{\alpha_0}} \left\langle \frac{Q}{\langle Q, U \rangle}, R \right\rangle.$$

Obviously,

$$\overset{*}{V} := \frac{Q}{\langle Q, U \rangle} \in \widetilde{\Theta}_x$$

and

$$\left\langle \overset{*}{V}, W \right\rangle \geq \alpha_0 \quad \text{for} \quad W \in \mathcal{A}_x.$$

So, $\overset{*}{V}$ is a maxmin strategy and $\alpha_0$ is the value of the game. Consider a decreasing sequence $(\alpha_k)$ with

$$\lim_{k \to \infty} \alpha_k = \alpha_0$$

and choose $\tilde{W}^k \in \mathcal{U}_{\alpha_k} \cap \widetilde{\mathcal{A}}_x$. By Lemma 3 there exists $W^k \in \mathcal{A}_x$ with

$$\tilde{W}^k \succeq W^k.$$

By Lemma 1 $(W^k)$ is bounded and so there exists a subsequence with

$$W := \lim_{k \in \mathbb{N}'} W^k.$$

Since $\mathcal{Q}_0$ is closed

$$\alpha_0 U - W = \lim_{k \in \mathbb{N}'} (\alpha_k U - W^k) \in \mathcal{Q}_0$$

i.e.

$$W \in \mathcal{U}_{\alpha_0}.$$

Now by Lemma 2 there exists $\overset{*}{W} \in \mathcal{A}_x$ with

$$W \succeq \overset{*}{W}$$

and therefore

$$\overset{*}{W} \in \mathcal{U}_{\alpha_0} \cap \mathcal{A}_x.$$

By definition of $\mathcal{U}_{\alpha_0}$

$$\alpha_0 \, U \succeq \overset{*}{W}$$

and it follows

$$\left\langle V, \overset{*}{W} \right\rangle \leq \alpha_0 \quad \forall \, V \in \tilde{\Theta}_x$$

and especially

$$\left\langle \overset{*}{V}, \overset{*}{W} \right\rangle = \alpha_0.$$

Then, $\left( \overset{*}{V}, \overset{*}{W} \right)$ is an equilibrium point of the linear game $(\widetilde{\Theta}_x, \mathcal{A}_x)$.

# References

BICKEL, P.J. and E.L. LEHMANN(1981). A minimax property of the sample mean in finite populations. Annals of Statistics 9, 1119-1122.

CHENG, C.S. and K.C. LI(1983). A minimax approach to sample surveys. Annals of Statistics 11,552-563.

GABLER, S.(1990). Minimax Solutions in Sampling from Finite Populations. Lecture Notes in Statistics 64, Springer: New York.

GABLER, S. and H. STENGER(2000). Minimax strategies in survey sampling. To appear in Journal of Statistical Planning and Inference.

HAJEK, J. (1959). Optimum strategy and other problems in probability sampling. Cas. Pest. Mat. 84, 387-473.

STENGER, H.(1990). Asymptotic minimaxity of the ratio strategy. Biometrika 77, 389-395.

STENGER, H. and S. GABLER(1996). A minimax property of Lahiri-Midzuno-Sen's sampling scheme. Metrika 43, 213-220.