

Sampling Space and Statistics on a Sample as an Element of a Semigroup

Žaneta Popeska

Faculty of Natural Sciences and Mathematics ,

Skopje, Republic of Macedonia

E-mail: zaneta@pmf.ukim.edu.mk

1. Sampling design

Let $B = \{b_1, b_2, \dots, b_N\}$ be a finite population, $S = S(B)$ or $S = \langle B \rangle$, a semigroup generated by B , and $U(B)$ a free semigroup generated by B . The elements of $U(B)$ will be denoted by $\sigma, \tau, \omega, \dots$, and the elements of a semigroup $S(B)$ by s, t, u, \dots . We define a **contents of** $s \in U$, by $C(\sigma) = \{b \mid b \in \sigma\}$ and a **length of** s by $L(\sigma) = n$ iff $\sigma = b_1 b_2 \dots b_n$, for $b_i \in B$.

If $S(B)$ is a semigroup generated by B , then there exists a unique homomorphism (which is epimorphism) $\psi: U(B) \rightarrow S(B)$ for which $\psi(b) = b$ for each $b \in B$. From now on we will use the symbol ψ only for this epimorphism. We define a **contents of** $s \in S$, by $C(s) = \{C(\sigma) \mid \psi(\sigma) = s\}$ and a **length of** s by $L(s) = \{n \mid n = L(\sigma), \psi(\sigma) = s\}$.

Let $p: S(B) \rightarrow [0, 1]$ be a real function.

Given $B, S(B)$, the triple $\mathbf{P} = (B, S(B), p)$ is called a sampling design if

- i) $\sum p(s) = 1$
- ii) $\forall b \in B, \exists s \in S(B)$, such that $b \in s$ and $p(s) > 0$.

The semigroup $S(B)$ is called a **sampling set**, and the real function p - a **design function**. The elements of S will be called S - **samples on** B . We say that the unit $b \in B$ belongs to a sample $s \in S$, and write $b \in s$ if s can be written as a product of elements of B in which b appears, i. e. $s = b_1 \cdot b_2 \cdot \dots \cdot b_n \in S(B)$, for $b_1, b_2, \dots, b_n \in B$ and $\exists i, 1 \leq i \leq n, b = b_i$. In other words $b \in s$ if and only if (iff) there is $\sigma \in U$ such that $b \in C(\sigma)$ and $\psi(\sigma) = s$.

2. Sampling space

Let B be a finite population and $\mathbf{Y}: B \rightarrow \mathfrak{R}$ a mapping. The vector $\mathbf{Y} = (Y_1, Y_2, \dots, Y_N)$, where $Y_i = \mathbf{Y}(b_i)$ is called a **population parameter**. Let U be the free semigroup and S a semigroup generated by B and $\psi: U \rightarrow S$ the epimorphism defined earlier. If we consider $\sigma \in U$ as a mapping $\sigma: b_i \rightarrow \sigma(b_i)$, let $\mathbf{y} = \mathbf{Y}\sigma = (Y_{\sigma(b_1)}, \dots, Y_{\sigma(b_n)})$. For $s \in S$ we define $\mathbf{Y}s = \{\mathbf{Y}\sigma \mid \sigma \in \psi^{-1}(\{s\})\}$ and $(s: \mathbf{Y}) = \{(\sigma, \mathbf{Y}\sigma) \mid \sigma \in \psi^{-1}(\{s\})\}$. Let $\Delta(S) = \{(s, \mathbf{y}) \mid s \in S_p, \mathbf{y} \in \mathfrak{R}^n, n \in L(s), \exists \mathbf{Y} \in \mathfrak{R}^N, \mathbf{y} \in \mathbf{Y}s\}$ and $\Delta^*(S) = \{(s: \mathbf{Y}) \mid s \in S, \mathbf{Y} \in \mathfrak{R}^N\}$. We define a relation \sim on $\Delta(U)$ by

$$(\sigma, \mathbf{y}) \sim (\tau, \mathbf{z}) \text{ iff } (\psi(\sigma) = \psi(\tau) \text{ and } \sigma_i = \tau_j \Rightarrow y_i = z_j).$$

This relation satisfies the following properties:

- 1° $(\sigma, \mathbf{y}) \sim (\tau, \mathbf{z}) \Rightarrow (\tau, \mathbf{z}) \sim (\sigma, \mathbf{y})$. ♦
- 2° $\forall (\sigma, \mathbf{y}) \in \Delta(U), (\sigma, \mathbf{y}) \sim (\sigma, \mathbf{y})$. ♦
- 3° If $\psi(\sigma) = \psi(\tau)$ and $\mathbf{y} = \mathbf{Y}\sigma, \mathbf{z} = \mathbf{Y}\tau$ for $\mathbf{Y} \in \mathfrak{R}^N$, then $(\sigma, \mathbf{y}) \sim (\tau, \mathbf{z})$. ♦
- 4° For any σ, τ for which $\psi(\sigma) = \psi(\tau)$, holds $(\sigma, \mathbf{Y}\sigma) \sim (\tau, \mathbf{Y}\tau) \forall \mathbf{Y} \in \mathfrak{R}^N$. ♦
- 5° The relation \sim does not to be transitive. ♦
- 6° If $\psi(\sigma) = \psi(\tau) \Rightarrow C(\sigma) = C(\tau)$ then the relation \sim is an equivalence relation. ♦

Let \approx be the transitive closure of \sim . Then \approx is an equivalence relation on Δ . The set $\Delta^\sim(S) = \Delta(U) / \approx$ is called a **data space** over the semigroup S and each equivalence class $(\sigma, \mathbf{y})^\sim$ is a **data** and will be denoted by $[\sigma, \mathbf{y}]$.

If for each $s \in S, \mathbf{y} \in \mathfrak{R}^n$, with $n \in L(s)$, for which there is $\sigma \in \psi^{-1}(s)$, such that $\ker \sigma \subseteq \ker \mathbf{y}$ we define a set $[s, \mathbf{y}] = \{[\tau, \mathbf{y}] \mid \tau \in \psi^{-1}(s)\}$ then $\Delta^\sim(S) = \bigcup_{s \in S} [s, \mathbf{y}]$.

- 7° For each $\sigma \in \psi^{-1}(s), (s, \mathbf{Y}s) \subseteq [\sigma, \mathbf{Y}\sigma]$. ♦

We say that a data $d=[\sigma,y] \in \Delta^{\sim}(S)$ is *consistent* with the parameter $Y \in \mathfrak{R}^N$ if and only if $(\sigma, Y\sigma) \in [\sigma, y]$. For each $Y \in \mathfrak{R}^N$ we define $p_Y^* : \Delta^*(S) \rightarrow [0,1]$ by

$$p_Y^*(s : Zs) = \begin{cases} p(s) & \text{if } Zs = Ys \\ 0 & \text{otherwise} \end{cases} .$$

8° The mapping p_Y^* is well defined and induces a probability measure on the algebra of subsets of $\Delta^*(S)$ defined by

$$P_Y^*(A) = \sum_{(s:Z) \in A} p_Y^*(s : Z) = \sum_{(s:Y) \in A} p(s) . \blacklozenge$$

Now we can define a mapping $p_Y : \Delta^{\sim}(S) \rightarrow [0,1]$ with

$$p_Y([\mathbf{s}, y]) = P_Y^* \{ (s : V) \mid (s : V) \subseteq [\mathbf{s}, y] \} .$$

$$9^\circ \quad p_Y([\mathbf{s}, y]) = \begin{cases} p(s) & (\mathbf{s}, y) \approx (\mathbf{s}, Y\mathbf{s}) \\ 0 & \text{otherwise} \end{cases} = \begin{cases} p(\mathbf{y}(\mathbf{s})) & [\mathbf{s}, y] = [\mathbf{s}, Y\mathbf{s}] \\ 0 & \text{otherwise} \end{cases} . \blacklozenge$$

10° The mapping p_Y induces a probability measure on the algebra of subsets of $\Delta^{\sim}(S)$, P_Y , defined for each $A \subseteq \Delta^{\sim}(S)$ by:

$$P_Y(A) = \sum_{[\mathbf{s}, y] \in A} p_Y([\mathbf{s}, y]) . \blacklozenge$$

For each $Y \in \mathfrak{R}^N$ we define subset $\Delta^{\sim}_Y \subseteq \Delta^{\sim}(S)$ with

$$\Delta^{\sim}_Y = \{ [\sigma, y] \mid [\sigma, y] \in \Delta^{\sim}(S), p_Y([\sigma, y]) > 0 \} .$$

11° The set Δ^{\sim}_Y is at most countable set. \blacklozenge

3. Statistics

Any mapping $F : \Delta^{\sim} \rightarrow \mathfrak{R}^m$ is a statistics over the design $\mathbf{P} = (B, \mathcal{S}(B), p)$. Now, having in mind that the set Δ^{\sim}_Y is at most countable and so is $F(\Delta^{\sim}_Y)$ one can define mathematical expectation of F , $E_Y(F)$ by

$$E_Y(F) = \sum_{w \in F(\Delta^{\sim}_Y)} w \cdot P_Y(F^{-1}(w))$$

Then

$$12^\circ \quad E_Y(F) = \sum_{\substack{s \in S_p \\ \text{for some} \\ t \in Y^{-1}(s)}} F([\mathbf{t}, Y\mathbf{t}]) \cdot p(s) . \blacklozenge$$

13° If F and G are statistics on $\mathbf{P} = (B, \mathcal{S}(B), p)$, for each $Y \in \mathfrak{R}^N$, and $\alpha, \beta \in \mathfrak{R}$ if the both sides exist,

$$E_Y(\alpha F + \beta G) = \alpha E_Y(F) + \beta E_Y(G) . \blacklozenge$$

REFERENCES

- Zaneta Popeska, 1998 Algebraic and combinatorial methods in sampling theory and experimental design, Ph.D. thesis, University St. Cyril & Metodij, Skopje, Macedonia
 Zaneta Popeska, 1994 Quotient sampling design, Macedonian academy of sciences & arts, XV 2, 37-48
 Cassel, Claes-Magnus, 1977 Foundations of inference in survey sampling", John & Sons, Inc.
 Jaroslav Hajek, 1981 Sampling from a finite population, Marcel Dekker