

A PROBABILITY DISTRIBUTION ASSOCIATED WITH EVENTS WITH MULTIPLE OCCURRENCES

John PANARETOS

University of Patras, Greece

Evdokia XEKALAKI

Athens University of Economics, Greece

Received August 1988

Revised October 1988

1. Introduction

In a recent paper Xekalaki and Panaretos (1989) introduce what they call the cluster negative binomial and cluster negative multinomial distributions in the context of cluster sampling. Clusters are represented by balls each marked with a number in the range $\{0, 1, 2, \dots, k\}$ where k is a positive integer. Balls are drawn one at a time and with replacement till a given number of zeroes is observed. The distribution of the total sum of the sampled numbers is defined to be the cluster negative binomial distribution. The case of straight sampling was considered by Panaretos and Xekalaki (1986a). Given a sample of n balls (clusters) numbered as previously and drawn with replacement the sum of the sampled numbers was shown to have a probability distribution termed the cluster binomial distribution. As mentioned by Xekalaki and Panaretos (1989), such urn schemes may be useful in modelling random mechanisms in the area of biology where the frequency distribution of cells forming clusters is of interest. Epidemiology is another area of possible application. There, the hypothesis of contagion is very common in determining the probability distribution of infected cases during an epidemic. In sampling-from-an-urn schemes this hypothesis can be accounted for by additional replacements which alter the structure of the urn. Panaretos and Xekalaki's (1986a) urn scheme giving rise to the stuttering generalized Waring distribution is pertinent to such a situation. There, balls marked $0, 1, 2, \dots, k$ are sampled with additional replacements till a given number of zeroes is observed. Of course among the various investigators who have used urn models for describing contagious biological data Polya (1931) seems to be a pioneer in the field.

This paper considers again sampling from an urn containing balls marked with an integer in the range $\{0, 1, 2, \dots, k\}$; $k > 0$, and studies the probability distribution of the sum of numbers shown on a sample of n balls drawn without replacement or with additional replacement(s). As a result, hypergeometric and

negative hypergeometric types of distributions are obtained. Their structural and distributional properties are studied and their association to other distributions arising under various schemes for sampling clusters are examined.

2. The cluster hypergeometric distribution

Consider an urn filled with α balls bearing number 0 and β_i balls bearing number i , $i = 1, 2, \dots, k$. A random sample of n balls is drawn without replacement. If X represents the sum of the numbers shown on the sampled balls then the following theorem can be shown.

Theorem 2.1. *The probability of the event $\{X = x\}$ is given by*

$$P(X = x) = \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, x_2, \dots, x_k, n - \sum_{i=1}^k x_i} \frac{\beta_1^{(x_1)} \dots \beta_k^{(x_k)} \alpha^{(n - \sum_{i=1}^k x_i)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}}, \tag{2.1}$$

$x = 0, 1, 2, \dots, nk$, where the symbol $\gamma^{(r)}$ stands for $\Gamma(\gamma + 1)/\Gamma(\gamma - r + 1)$, $r = 0, 1, \dots, \gamma$, $\gamma > 0$.

Proof. Let X_i denote the number of balls bearing number i in the sample of n balls. Then obviously $\{X = x\}$ is the union of all the outcomes $\{X_1 = x_1, \dots, X_k = x_k\}$ for which $\sum_{i=1}^k ix_i = x$. i.e.

$$\{X = x\} = \bigcup_{\sum_{i=1}^k ix_i = x} \{X_1 = x_1, \dots, X_k = x_k\}.$$

Then (2.1) follows as a consequence of the fact that the vector (X_1, X_2, \dots, X_k) follows the multivariate hypergeometric distribution with parameters $n, \alpha, \beta_1, \dots, \beta_k$. \square

Theorem 2.2. *The function defined by (2.1) is a proper probability function.*

Proof. It is sufficient to show that $\sum_{x=0}^{nk} P(X = x) = 1$. Indeed

$$\sum_{x=0}^{nk} P(X = x) = \sum_{x=0}^{nk} \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, x_2, \dots, x_k, n - \sum_{i=1}^k x_i} \frac{\beta_1^{(x_1)} \dots \beta_k^{(x_k)} \alpha^{(n - \sum_{i=1}^k x_i)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}}.$$

Letting $x_i \rightarrow x_i$ and $x \rightarrow x + \sum_{i=1}^k (i - 1)x_i$ we obtain

$$\sum_{x=0}^n P(X = x) = \sum_{x=0}^n \binom{n}{x} \frac{\alpha^{(n-x)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}} \sum_{\sum_{i=1}^k x_i = x} \binom{x}{x_1, \dots, x_k} \beta_1^{(x_1)} \dots \beta_k^{(x_k)} = 1.$$

Hence the theorem has been established. \square

Note. One can easily observe that for $k = 1$ the probability distribution defined by (2.1) reduces to the ordinary hypergeometric distribution. Moreover, by the definition of the descending factorial the definition of (2.1) can be extended to include the case of non integer values for the parameters α and β_i , $i = 1, 2, \dots, k$.

Definition 2.1. Let n, k be fixed positive integers and let $\alpha, \beta_i, i = 1, 2, \dots, k$ be positive real numbers. A random variable X will be said to have the cluster hypergeometric distribution with parameters $k, n, \alpha, \beta_1, \dots, \beta_k$ ($X \sim H_k(n; \alpha, \beta_1, \dots, \beta_k)$) if its probability function is given by (2.1).

Steyn's (1956) univariate factorial multinomial distribution of the number of occurrences of an event E when sampling from a population whose items cause multiple occurrences of E can also be obtained from (2.1) if the proportion of balls marked i is allowed to be an arbitrary real number in $(0, 1)$ instead of a rational number in $(0, 1)$.

Theorem 2.3. Let X be a non-negative integer-valued random variable having the cluster hypergeometric distribution with probability function as given by (2.1). Then

$$(i) \quad G_X(s) = \frac{\alpha^{(n)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}} F_D(-n; -\beta_1, -\beta_2, \dots, -\beta_k; \alpha - n + 1; t, t^2, \dots, t^k) \quad (2.2)$$

where $G_X(s)$ denotes the probability generating function of X , and F_D is Lauricella's generalized hypergeometric function derived by

$$F_D(a; b_1, \dots, b_k; c; z_1, \dots, z_k) = \sum_{r=0}^{\infty} \sum_{\sum_{i=1}^k r_i=r} \frac{a_{(\sum_{i=1}^k r_i)} (b_1)_{(r_1)} \dots (b_k)_{(r_k)} z_1^{r_1} \dots z_k^{r_k}}{c_{(\sum_{i=1}^k r_i)}} \quad (2.3)$$

(here $a_{(r)}$ stands for $\Gamma(a+r)/\Gamma(a)$, note that when $a < 0$ the series is terminating, i.e. r runs in the range $0, 1, \dots, [a]$);

$$(ii) \quad E(X) = \frac{n}{\alpha + \sum_{i=1}^k \beta_i} \sum_{i=1}^k i\beta_i; \quad (2.4)$$

$$(iii) \quad V(X) = \frac{n(\alpha + \sum_{i=1}^k \beta_i - n) \{ \sum_{i=1}^k i^2 \beta_i (\alpha + \sum_{j \neq i}^k \beta_j) - 2 \sum_{i < j} ij \beta_i \beta_j \}}{(\alpha + \sum_{i=1}^k \beta_i)^2 (\alpha + \sum_{i=1}^k \beta_i - 1)}. \quad (2.5)$$

Proof. (i) Applying the transformation $x_i \rightarrow x_i, x \rightarrow x + \sum_{i=1}^k (i-1)x_i$ we have

$$\begin{aligned} G_X(s) &= \sum_{x=0}^n \sum_{\sum_{i=1}^k x_i=x} \binom{n}{x_1, \dots, x_k, n-x} \frac{\beta_1^{(x_1)} \dots \beta_k^{(x_k)} \alpha^{(n-x)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}} s^{\sum_{i=1}^k i x_i} \\ &= \frac{\alpha^{(n)}}{(\alpha + \sum_{i=1}^k \beta_i)^{(n)}} \sum_{x=0}^n \sum_{\sum_{i=1}^k x_i=x} \frac{\beta_1^{(x_1)} \dots \beta_k^{(x_k)} (\alpha - n)^{(-\sum_{i=1}^k x_i)}}{(1-n)^{(-\sum_{i=1}^k x_i)}} \prod_{i=1}^k \frac{s^{i x_i}}{x_i!} \end{aligned}$$

which is equivalent to (2.2).

The proofs of parts (ii) and (iii) follow as a result of the fact that (X_1, X_2, \dots, X_k) has the multivariate hypergeometric distribution with $E(X_i)$ and $E(X_i X_j)$ given by the appropriate well-known formulas for $i = 1, 2, \dots, k; j = 1, 2, \dots, k$. \square

As is well known in the context of urn models the ordinary hypergeometric distribution is the analogue of the binomial distribution when sampling is done without replacement and the former has the latter as a limiting form when the numbers of balls are increased so that their proportion tend to some value $p \in (0, 1)$. It was mentioned in the introduction that a scheme analogous to that considered in this section but whereby sampling is done with replacement gives rise to Panaretos and Xekalaki's (1986b) cluster

binomial distribution. This is the distribution defined by the probability function

$$P(X = x) = \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, x_2, \dots, x_k, n - \sum_{i=1}^k x_i} q^{n - \sum_{i=1}^k x_i} \prod_{i=1}^k p_i^{x_i}, \tag{2.6}$$

$$q = 1 - \sum_{i=1}^k p_i, \quad 0 < p_i < 1, \quad i = 1, 2, \dots, k, \quad x = 0, 1, \dots, nk.$$

It would therefore be interesting to examine whether the distribution in (2.6) can serve as an approximation of (2.1). The following theorem shows that this is indeed the case.

Theorem 2.4. *Let X, Y be two random variables whose probability functions are given by (2.1) and (2.6) respectively. If $\alpha \rightarrow \infty, \beta_i \rightarrow \infty, i = 1, 2, \dots, k$ so that $\beta_i / (\alpha + \sum_{i=1}^k \beta_i) \rightarrow p_i$ where p_i are fixed values in the range $(0, 1), i = 1, 2, \dots, k$, then the distribution of X tends to the distribution of Y .*

Proof. Let \lim_H denote limit as $\alpha \rightarrow \infty, \beta_i \rightarrow \infty$ so that $\beta_i / (\alpha + \sum_{i=1}^k \beta_i) \rightarrow p_i, i = 1, 2, \dots, k$. Then

$$\begin{aligned} \lim_H P(X = x) &= \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, \dots, x_k, n - \sum_{i=1}^k x_i} \\ &\quad \times \lim_H \left\{ \left(\frac{\alpha}{\alpha + \sum_{i=1}^k \beta_i} \right)^{n - \sum_{i=1}^k x_i} \prod_{i=1}^k \left(\frac{\beta_i}{\alpha + \sum_{i=1}^k \beta_i} \right)^{x_i} \right\} \\ &= \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, \dots, x_k, n - \sum_{i=1}^k x_i} \left(1 - \sum_{i=1}^k p_i \right)^{n - \sum_{i=1}^k x_i} \prod_{i=1}^k p_i^{x_i}. \end{aligned}$$

Hence the result. \square

3. The cluster negative hypergeometric distribution

Suppose now that we wish to sample n balls one at a time from the urn of the previous section, altering its composition at each stage by adding a ball of the same type as the sampled one before the next draw is made. As far as applications are concerned this allows for the hypothesis of contagion since each added ball increases the probability of sampling a ball of the same type at the next draw.

So we start off with an urn containing α balls marked 0 and β_i balls marked $i, i = 1, 2, \dots, k$. A ball is drawn, its number is recorded and the ball is returned to the urn with one additional ball of the same type. Let X be the sum of the numbers observed in a sample of n balls sampled in the described manner. Then the following theorems can be shown following arguments similar to those used to prove Theorems 2.1 and 2.2.

Theorem 3.1. *The probability of the event $\{X = x\}$ is given by*

$$P(X = x) = \sum_{\sum_{i=1}^k ix_i = x} \binom{n}{x_1, \dots, x_k, x - \sum_{i=1}^k x_i} \frac{(\beta_1)_{(x_1)} \cdots (\beta_k)_{(x_k)} \alpha_{(n - \sum_{i=1}^k x_i)}}{(\alpha + \sum_{i=1}^k \beta_i)_{(n)}}, \tag{3.1}$$

$$x = 0, 1, 2, \dots, nk. \quad \square$$

Theorem 3.2. *Let X be a random variable defined as in Theorem 3.1. Then $\sum_{x=0}^{nk} P(X = x) = 1$. \square*

The result of this theorem implies that (3.1) defines a probability distribution which for $k = 1$ coincides with the ordinary negative hypergeometric distribution.

Definition 3.1. Let n, k be fixed positive integers and let $\alpha, \beta_i, i = 1, 2, \dots, k$, be positive real numbers. A random variable X will be said to have the cluster negative hypergeometric distribution with parameters $k, n, \alpha, \beta_1, \dots, \beta_k$ ($X \sim H_k^-(x; n, \alpha, \beta_1, \dots, \beta_k)$) if its probability function is given by (3.1).

As before, if instead of the rational proportion $\beta_i/(\alpha + \sum_{i=1}^k \beta_i)$ of balls in the urn we start off with a proportion equal to any real number p_i in $(0, 1)$, the urn scheme considered in this section provides an alternative derivation of Steyn's (1956) univariate negative factorial multinomial distribution of occurrences of an event caused repeatedly by inversely sampled items of a population.

Theorem 3.3. Let X be a non-negative integer-valued random variable with probability function as given by (3.1). Then

$$(i) \quad G_X(s) = \frac{\alpha^{(n)}}{(\alpha + \sum_{i=1}^k \beta_i)_{(n)}} F_D(-n; \beta_1, \beta_2, \dots, \beta_k; 1 - \alpha - n; s, s^2, \dots, s^k), \quad (3.2)$$

$$(ii) \quad E(X) = \frac{n}{\alpha + \sum_{i=1}^k \beta_i} \sum_{i=1}^k i \beta_i, \quad (3.3)$$

$$(iii) \quad V(X) = \frac{n(\alpha + \sum_{i=1}^k \beta_i + n) \{ -2 \sum_{i < j} i j \beta_i \beta_j + \sum_{i=1}^k i^2 \beta_i (\alpha + \sum_{j \neq i} \beta_j) \}}{(\alpha + \sum_{i=1}^k \beta_i)^2 (\alpha + \sum_{i=1}^k \beta_i + 1)}. \quad (3.4)$$

Proof. The proof is analogous to that of Theorem 2.3. \square

The two theorems that follow demonstrate the fact that the cluster negative hypergeometric distribution and the cluster binomial distribution relate to each other in the same manner as their ordinary forms (for $k = 1$) do. In particular, it is shown that (3.1) arises as a Dirichlet mixture of (2.6). Moreover (2.6) is shown to provide an approximation to (3.1) for certain limiting values of its parameters.

Theorem 3.4. Let X be a non-negative integer valued random variable having the cluster binomial distribution with probability function as given by (2.6). Suppose that the parameter vector $\mathbf{p} = (p_1, p_2, \dots, p_k)$ has itself a distribution with probability density function of the form

$$f(\mathbf{p}) = \frac{\Gamma(\sum_{i=1}^k \beta_i + \alpha)}{\Gamma(\alpha) \prod_{i=1}^k \Gamma(\beta_i)} p_1^{\beta_1 - 1} \dots p_k^{\beta_k - 1} \left(1 - \sum_{i=1}^k p_i \right)^{\alpha - 1}, \quad (3.5)$$

$\alpha, p_i, \beta_i > 0, i = 1, 2, \dots, k, \sum_{i=1}^k p_i < 1$. Then the resulting distribution of X is the cluster negative hypergeometric distribution with probability function given by (3.1).

Proof.

$$\begin{aligned} P(X = x) &= \int_0^1 \dots \int_0^1 P(X = x | \mathbf{p}) f(\mathbf{p}) \, dp_1 \dots dp_k \\ &= \frac{\Gamma(\sum_{i=1}^k \beta_i + \alpha)}{\Gamma(\alpha) \prod_{i=1}^k \Gamma(\beta_i)} \int_0^1 \dots \int_0^1 \sum_{\sum_{i=1}^k i x_i = x} \binom{n}{x_1, \dots, x_k, n - \sum_{i=1}^k x_i} \\ &\quad \times \left(1 - \sum_{i=1}^k p_i \right)^{n + \alpha - \sum_{i=1}^k x_i - 1} \left\{ \prod_{i=1}^k p_i^{x_i + \beta_i - 1} \right\} dp_1 \dots dp_k \end{aligned}$$

$$= \frac{\Gamma(\sum_{i=1}^k \beta_i + \alpha)}{\Gamma(\alpha) \prod_{i=1}^k \Gamma(\beta_i)} \sum_{\sum_{i=1}^k x_i = x} \binom{n}{x_1, \dots, x_k, n - \sum_{i=1}^k x_i} \\ \times \frac{\{\prod_{i=1}^k \Gamma(x_i + \beta_i)\} \Gamma(n + \alpha - \sum_{i=1}^k x_i)}{\Gamma(\sum_{i=1}^k \beta_i + n + \alpha)}.$$

Hence, the result. \square

Theorem 3.5. Let X be a random variable distributed according to the cluster negative hypergeometric distribution with probability function as given by (3.1) and let Y be another random variable having the cluster binomial distribution with probability function given by (2.6). Then, if $\alpha \rightarrow \infty$ and $\beta_i \rightarrow \infty$ so that $\beta_i/(\alpha + \sum_{i=1}^k \beta_i) \rightarrow p_i, 0 < p_i < 1, i = 1, 2, \dots, k$, the distribution of X tends to the distribution of Y .

Proof. Let \lim_H stand for limit as $\alpha \rightarrow \infty, \beta_i \rightarrow \infty$ so that $\beta_i/(\alpha + \sum_{i=1}^k \beta_i) \rightarrow p_i, i = 1, 2, \dots, k$. Then

$$\lim_H P(X = x) = P(Y = x)$$

since

$$\lim_H \frac{\alpha_{(n - \sum_{i=1}^k x_i)} \prod_{i=1}^k (\beta_i)_{(x_i)}}{(\alpha + \sum_{i=1}^k \beta_i)_{(n)}} = \lim_H \left(\frac{\alpha}{\alpha + \sum_{i=1}^k \beta_i} \right)^{n - \sum_{i=1}^k x_i} \prod_{i=1}^k \lim_H \left(\frac{\beta_i}{\alpha + \sum_{i=1}^k \beta_i} \right)^{x_i} \\ = \left(1 - \sum_{i=1}^k p_i \right)^{n - \sum_{i=1}^k x_i} \prod_{i=1}^k p_i^{x_i}.$$

Hence the theorem has been established. \square

The sampling scheme considered in this section can be slightly altered so as to yield a more general type of distribution. In particular, if each sampled ball is returned to the urn along with c additional balls of the same type ($c \geq 1$) before the next ball is drawn then X , the sum of the n sampled numbers follows a distribution with a probability function given by

$$P(X = x) = \sum_{\sum_{i=1}^k x_i = x} \binom{n}{x_1, \dots, x_k, n - \sum_{i=1}^k x_i} \frac{(\alpha/c)_{(n - \sum_{i=1}^k x_i)} \prod_{i=1}^k (\beta_i/c)_{(x_i)}}{((\sum_{i=1}^k \beta_i + \alpha)/c)_{(n)}}, \tag{3.6}$$

$\alpha, c, \beta_i > 0, i = 1, 2, \dots, k, x = 0, 1, 2, \dots, nk$. For $k = 1$ this reduces to the ordinary Polya distribution.

Definition 3.2. A non-negative, integer-valued random variable X taking values in $\{0, 1, \dots, nk\}$, where k is a positive integer is said to have the cluster Polya distribution with parameters c, α and $\beta_i, i = 1, 2, \dots, k$, if its probability function is given by (3.6).

Note that when $c = 0$, (3.6) reduces to the cluster binomial as this is defined by (2.6) for $p_i = \beta_i/(\alpha + \sum_{i=1}^k \beta_i), i = 1, 2, \dots, k$. Moreover when $c = -1$, (3.6) reduces to the cluster hypergeometric with probability function as given by (2.1).

Hence (3.6) unifies the derivation of the distributions that arise in the context of direct cluster sampling, namely the cluster binomial of Panaretos and Xekalaki (1986a) and the cluster hypergeometric and the cluster negative hypergeometric distributions introduced in the present paper.

Acknowledgement

Part of this work was done while the authors were visiting the Institute of Statistical Mathematics in Tokyo as visiting fellows. The financial support provided by the Institute as well as the interesting conversations on the subject with Drs. Shimizu, Hirano and Aki on the topic are gratefully acknowledged.

References

- Panaretos, J. and E. Xekalaki (1986a), The stuttering generalized Waring distribution, *Statist. Probab. Lett.* **4**, 313–318.
- Panaretos, J. and E. Xekalaki (1986b), On generalized binomial and multinomial distributions and their relation to generalized Poisson distributions, *Ann. Inst. Statist. Math. A* **38**, 223–231.
- Polya, G. (1931), Sur quelques points de la theorie des probaliteés, *Ann. Inst. H. Poincaré* **1**, 117–161.
- Steyn, H.S. (1956), On the univariable series $F(t) = F(a; b_1, b_2, \dots, b_k; c; t, t^2, \dots, t^k)$ and its application in probability theory, *Proc. Kon. Ned. Akad. Wetensch. Ser. A* **59**, 190–197.
- Xekalaki, E. and J. Panaretos (1989), On some distributions arising in inverse cluster sampling, *Comm. Statist. A—Theory Methods* **18** (1) 355–366.