

SOME ASYMPTOTIC RESULTS FOR OCCUPANCY PROBLEMS¹

BY LARS HOLST

University of Wisconsin, Madison

Suppose n balls are placed into N cells with arbitrary probabilities. Limit distributions for the number of empty cells are considered when $N \rightarrow \infty$ and $n \rightarrow \infty$ in such a way that $n/N \rightarrow \infty$. Limit distributions for the number of balls to achieve exactly b empty cells are obtained for $N \rightarrow \infty$ with b fixed and for $b \rightarrow \infty$ with $b/N \rightarrow 0$.

1. Introduction. Suppose that balls are thrown independently of each other into N cells so that each ball has probability p_k of falling into the k th cell, $p_1 + \dots + p_N = 1$. Let Y_n denote the number of empty cells after n throws and let T_b denote the throw on which for the first time exactly b cells remain empty, $0 \leq b < N$. The symmetrical case $p_1 = \dots = p_N = 1/N$ is discussed for example, in Feller (1968) under occupancy or waiting time problems. For an expository paper on these and related problems, see Kolchin and Chistyakov (1974).

Depending on how $b, n, N \rightarrow \infty$, different asymptotic distributions for Y_n and T_b can be obtained; see, for example, Holst (1971) and, for the symmetric case, Samuel-Cahn (1974). In this paper some remaining problems are investigated for the nonsymmetrical case.

To give precise meanings for the limits obtained, double sequences $(p_{kN})_N$, $(Y_{nN})_N$ are considered. But in order to simplify the notation the extra index N will usually be omitted.

2. A bounded number of empty cells. The following limit theorem for Y_n , the number of empty cells after n throws, was proved by Sevastyanov (1972).

THEOREM 1. *If the p 's are such that*

$$(2.1) \quad \max_{1 \leq k \leq N} (1 - p_k)^n \rightarrow 0$$

and

$$(2.2) \quad E(Y_n) = \sum_{k=1}^N (1 - p_k)^n \rightarrow m < \infty,$$

then

$$(2.3) \quad P(Y_n = y) \rightarrow m^y \cdot e^{-m}/y!,$$

or equivalently

$$(2.4) \quad Y_n \Rightarrow P o(m), \quad \text{when } N \rightarrow \infty.$$

Received April 8, 1976; revised February 24, 1977.

¹ Sponsored by the United States Army under Contract No. DAAG29-75-C-0024.

AMS 1970 subject classifications. Primary 60F05; Secondary 60C05.

Key words and phrases. Occupancy problems, coupon collectors problem, limit theorems.

REMARK. When the p 's are equal an expression for $P(Y_n = y)$ can be obtained from which (2.3) can be derived by elementary methods; see Feller (1968). In this case (2.1) and (2.2) are replaced by

$$(2.5) \quad N \cdot \exp(-n/N) \rightarrow m < \infty .$$

For T_b , the number of balls until b empty cells remain, we have:

THEOREM 2. If b is a fixed integer and for some fixed numbers C and D

$$(2.6) \quad 0 < C \leq Np_k \leq D < \infty , \quad \text{for all } k \text{ and } N ,$$

then, when $N \rightarrow \infty$,

$$(2.7) \quad \sum_{k=1}^N (1 - p_k)^{T_b} \Rightarrow \frac{1}{2}\chi^2(2(b + 1)) ,$$

and

$$(2.8) \quad \sum_{k=1}^N \exp(-T_b p_k) \Rightarrow \frac{1}{2}\chi^2(2(b + 1)) .$$

Before we prove the theorem let us consider the functions

$$(2.9) \quad f(t) = f_N(t) = \sum_{k=1}^N (1 - p_k)^t , \quad t > 0 ,$$

and

$$(2.10) \quad g(t) = g_N(t) = \sum_{k=1}^N \exp(-tp_k) .$$

LEMMA 1. If condition (2.6) is satisfied, $y > 0$ is a fixed number, and $t = t_N = t(y)$ is defined by the equation

$$(2.11) \quad f(t) = y ,$$

then

$$(2.12) \quad 0 < C \leq \liminf_{N \rightarrow \infty} N(\log N)/t_N \leq \limsup_{N \rightarrow \infty} N(\log N)/t_N \leq D < \infty ,$$

and when $N \rightarrow \infty$

$$(2.13) \quad f([t]) \rightarrow y ,$$

$$(2.14) \quad \max_{1 \leq k \leq N} (1 - p_k)^{[t]} \rightarrow 0 ,$$

$$(2.15) \quad g(t) \quad \text{and} \quad g([t]) \rightarrow y ,$$

where $[t]$ denotes the integer part of t .

LEMMA 2. If f is replaced by g and g by f in Lemma 1, then the same conclusions hold.

PROOF OF LEMMA 1. From condition (2.6), it follows that

$$(2.16) \quad y = \sum_{k=1}^N (1 - p_k)^t \geq N \cdot (1 - D/N)^t .$$

Hence for $\epsilon > 0$ and N sufficiently large,

$$(2.17) \quad \log y \geq \log N - t \cdot (D + \epsilon)/N ,$$

and therefore

$$(2.18) \quad \begin{aligned} D + \varepsilon &= (D + \varepsilon) \lim_{N \rightarrow \infty} (1/(1 - \log y/\log N)) \\ &\geq \limsup_{N \rightarrow \infty} N \log N/t_N, \end{aligned}$$

which proves the right inequality of (2.12).

To prove the left inequality of (2.12), note that by (2.6)

$$(2.19) \quad \log y \leq \log N - t \log(1 - C/N) \leq \log N - tC/N.$$

From this it follows that

$$(2.20) \quad C = C \lim_{N \rightarrow \infty} (1 - \log y/\log N)^{-1} \leq \liminf_{N \rightarrow \infty} N \log N/t_N.$$

Using (2.6) and (2.11) we get

$$(2.21) \quad (1 - D/N)^{-1}y \geq f([t]) \geq y,$$

which proves (2.13).

Combining (2.6) and (2.12) shows that, for some $K > 0$ and N sufficiently large,

$$(2.22) \quad \max (1 - p_k)^{[t]} \leq (1 - C/N)^{[t]} \leq (1 - C/N)^{KN \log N} \rightarrow 0, \quad N \rightarrow \infty,$$

proving (2.14).

From (2.6) and (2.12) it follows that for some constant K

$$(2.23) \quad |1 - e^{-tp_k}/(1 - p_k)^t| \leq K \cdot \log N/N,$$

and therefore

$$(2.24) \quad \begin{aligned} |f(t) - g(t)| &\leq \sum_1^N (1 - p_k)^t \cdot |1 - e^{-tp_k}/(1 - p_k)^t| \\ &\leq K \sum_1^N (1 - p_k)^t (\log N)/N = Ky(\log N)/N \rightarrow 0, \end{aligned}$$

which proves (2.15).

PROOF OF LEMMA 2. The proof is essentially the same as that for Lemma 1.

PROOF OF THEOREM 2. From the definitions it follows that

$$(2.25) \quad Y_n \leq b \Leftrightarrow T_b \leq n,$$

and therefore

$$(2.26) \quad P(Y_n \leq b) = P(T_b \leq n) = P(f(T_b) \geq f(n)).$$

Let $y > 0$ be fixed and define $n = [t]$ with $t = t(y)$ as in Lemma 1. According to Lemma 1 the assumptions of Theorem 1 are satisfied. Hence

$$(2.27) \quad P(f(T_b) \geq y) = P(Y_n \leq b) \rightarrow P(Y \leq b),$$

where Y is $P o(y)$. Furthermore it is well-known that

$$(2.28) \quad P(Y \leq b) = P(\frac{1}{2}\chi^2(2(b+1)) \geq y);$$

(2.27) and (2.28) prove (2.7). From Lemma 2 the assertion (2.8) follows.

REMARK. When the p 's are equal the theorem can be written

$$(2.29) \quad N \cdot (1 - 1/N)^{T_b} \Rightarrow \frac{1}{2}\chi^2(2(b + 1)),$$

and therefore

$$(2.30) \quad T_b/N - \log N \Rightarrow \log(\frac{1}{2}\chi^2(2(b + 1))).$$

This result was found by Baum and Billingsley (1965) using complicated calculations. From the result in Feller (1968) and the method of proof of Theorem 2, (2.29) and (2.30) follow. A consequence of (2.30) is

$$(2.31) \quad T_b/N \log N \rightarrow 1 \quad \text{in probability as } N \rightarrow \infty.$$

Now (2.31) will be generalized. First introduce the distribution function

$$(2.32) \quad H_N(x) = \#(p_k; Np_k \leq x)/N.$$

LEMMA 3. If $t = t_N = t(y)$ is defined by

$$(2.33) \quad g(t) = g_N(t_N) = y > 0,$$

and there exists a distribution function $H(x)$ on $[C, D]$ such that

$$(2.34) \quad H_N(x) \rightarrow H(x), \quad N \rightarrow \infty,$$

and

$$(2.35) \quad 0 < C = \inf \{x; H(x) > 0\},$$

then for $1/C > \varepsilon > 0$,

$$(2.36) \quad g_N((\varepsilon + 1/C)(N \log N)) \rightarrow 0,$$

and

$$(2.37) \quad g_N((-\varepsilon + 1/C)(N \log N)) \rightarrow +\infty$$

as $N \rightarrow \infty$.

PROOF. From the definitions it follows that

$$(2.38) \quad \begin{aligned} 0 < y = g_N(t_N) &= N \cdot \int_C^D \exp(-t_N x/N) dH_N(x) \\ &= \int_C^D \exp((1 - t_N x/N \log N) \log N) dH_N(x). \end{aligned}$$

Consider

$$(2.39) \quad g_N((\varepsilon + 1/C)N \log N) = \int_C^D \exp((1 - x(1 + \varepsilon C)/C) \log N) dH_N(x).$$

Now, for $C \leq x \leq D$ it is true that $1 - x(1 + \varepsilon C)/C < 0$ and therefore the exponent in (2.39) is negative, so the integral tends to 0 when $N \rightarrow \infty$, which proves (2.36).

In a similar way (2.38) can be proved.

COROLLARY (to Theorem 2). If the conditions (2.34) and (2.35) are satisfied, then

$$(2.40) \quad T_b/N \log N \rightarrow 1/C \quad \text{in probability as } N \rightarrow \infty.$$

PROOF. Let $\varepsilon_1 > 0$ and $\varepsilon_2 > 0$ be given. Take a $\delta > 0$ so that

$$(2.41) \quad P(\frac{1}{2}\chi^2(2(b + 1)) < \delta) < \varepsilon_2/2 .$$

For N sufficiently large it follows from Theorem 2 that

$$(2.42) \quad P(g_N(T_b) < \delta) < \varepsilon_2/2$$

and from Lemma 3 that

$$(2.43) \quad g_N((\varepsilon_1 + 1/C)(N \log N)) < \delta .$$

Hence

$$(2.44) \quad P(T_b/N \log N > \varepsilon_1 + 1/C) = P(g_N(T_b) < g_N((\varepsilon_1 + 1/C)(N \log N))) \\ \leq P(g_N(T_b) < \delta) < \varepsilon_2/2 .$$

In a similar way we prove

$$(2.45) \quad P(T_b/N \log N < -\varepsilon_1 + 1/C) < \varepsilon_2/2 .$$

By (2.44) and (2.45) the assertion follows.

3. A small fraction of empty cells. As above, Y_n denotes the number of empty cells after n throws.

THEOREM 3. *If*

$$(3.1) \quad 0 < C \leq Np_k \leq D < \infty \quad \text{for all } k \text{ and } N ,$$

$$(3.2) \quad n/N \rightarrow \infty ,$$

$$(3.3) \quad f(n) = E(Y_n) = \sum_{k=1}^N (1 - p_k)^n \rightarrow +\infty ,$$

then when $n \rightarrow \infty$,

$$(3.4) \quad (Y_n - f(n))/(f(n))^{\frac{1}{2}} \Rightarrow N(0, 1) ,$$

and

$$(3.5) \quad (Y_n - g(n))/(g(n))^{\frac{1}{2}} \Rightarrow N(0, 1) ,$$

where

$$(3.6) \quad g(n) = \sum_{k=1}^N \exp(-np_k) .$$

PROOF. From (3.1) and (3.3) it follows that

$$(3.7) \quad \sum_{i=1}^N (1 - p_k)^n \leq N \cdot (1 - C/N)^n \rightarrow +\infty ;$$

hence

$$(3.8) \quad n/N \log N = O(1) .$$

By (3.1), (3.2), and (3.8) there exists a constant K such that

$$(3.9) \quad |f(n) - g(n)| \leq K \cdot (n/N) \cdot \exp(-Cn/N) \rightarrow 0 .$$

Hence it is sufficient to prove (3.5). This will be done using a technique similar to that of Karlin (1967).

Let $\{X(t)\}$ be a Poisson process with unit parameter and assume that at each event of the Poisson process we independently choose one of the N cells according to $\{p_k\}_{k=1}^N$ and place a ball in it. For $1 \leq k \leq N$ let

$$(3.10) \quad W_k(t) = \text{number of balls in cell } k \text{ at time } t.$$

It is well known that for each t , $\{W_k(t)\}$ are independent random variables with $W_k(t)$ distributed as $Po(p_k t)$. Let

$$(3.11) \quad \begin{aligned} I(y) &= 1 && \text{if } y = 0 \\ &= 0 && \text{otherwise.} \end{aligned}$$

Define

$$(3.12) \quad Y(t) = \sum_{k=1}^N I(W_k(t)).$$

Then $Y(t)$ is the number of empty cells at time t , and by (3.3) and (3.9)

$$(3.13) \quad E(Y(n)) = g(n) = \sum_{k=1}^N e^{-np_k} \rightarrow \infty.$$

Thus, by the central limit theorem and the independence of the $\{W_k(t)\}$,

$$(3.14) \quad (Y(n) - g(n))/g(n)^{\frac{1}{2}} \Rightarrow N(0, 1).$$

We now need to show for all x

$$(3.15) \quad |P((Y(n) - g(n))/g(n)^{\frac{1}{2}} \leq x) - P((Y_n - g(n))/g(n)^{\frac{1}{2}} \leq x)| \rightarrow 0$$

as $n \rightarrow \infty$. But

$$(3.16) \quad P((Y(n) - g(n))/g(n)^{\frac{1}{2}} \leq x) = \sum_{j=0}^{\infty} P((Y_j - g(n))/g(n)^{\frac{1}{2}} \leq x) e^{-n} \cdot n^j/j!.$$

Let $\delta > 0$. Then by the central limit theorem there exists a $A > 0$ such that for all x and n sufficiently large,

$$(3.17) \quad \begin{aligned} |P((Y(n) - g(n))/g(n)^{\frac{1}{2}} \leq x) \\ - \sum_{|j-n| \leq An^{\frac{1}{2}}} P((Y_j - g(n))/g(n)^{\frac{1}{2}} \leq x) e^{-n} \cdot n^j/j!| < \delta. \end{aligned}$$

Let $\epsilon > 0$ and suppose we can prove that

$$(3.18) \quad \sup_{|j-n| < An^{\frac{1}{2}}} P(|Y_n - Y_j| > \epsilon(g(n))^{\frac{1}{2}}) = o(1).$$

It then will follow from (3.17) and (3.18) that

$$(3.19) \quad \begin{aligned} P((Y(n) - g(n))/g(n)^{\frac{1}{2}} < x - \epsilon) &< -\delta + o(1) \\ &\leq P((Y_n - g(n))/g(n)^{\frac{1}{2}} \leq x) \\ &\leq P((Y(n) - g(n))/g(n)^{\frac{1}{2}} \leq x + \epsilon) + \delta + o(1), \end{aligned}$$

which in turn will establish the theorem by the continuity of the normal distribution. So it remains only to prove (3.18).

Markov's inequality yields

$$(3.20) \quad P(|Y_n - Y_j| > \epsilon(g(n))^{\frac{1}{2}}) \leq E|Y_n - Y_j|/\epsilon(g(n))^{\frac{1}{2}}.$$

Assume first that $j > n$, so $j = n + i$, $0 < i < An^{\frac{1}{2}}$. Since $Y_n \geq Y_j$,

$$(3.21) \quad \begin{aligned} E|Y_n - Y_{n+i}| &= E(Y_n) - E(Y_{n+i}) = \sum_{k=1}^N (e^{-np_k} - e^{-(n+i)p_k}) \\ &\leq \sum_{k=1}^N e^{-np_k}(1 - \exp(-An^{\frac{1}{2}}p_k)). \end{aligned}$$

Since $\max_k p_k \leq D/N$ and $n/N \log N = O(1)$, $n^{1/2} \max p_k \rightarrow 0$ and so for some constant D_2 ,

$$(3.22) \quad \sup_{0 \leq i \leq An^{1/2}} E|Y_n - Y_{n+i}| \leq D_2 g(n) n^{1/2} / N.$$

Similarly

$$(3.23) \quad \sup_{0 \leq i \leq An^{1/2}} E|Y_n - Y_{n-i}| \leq D_2 g(n) n^{1/2} / N.$$

Thus

$$(3.24) \quad \sup_{|k-n| < An^{1/2}} P(|Y_n - Y_k| > \varepsilon(g(n))^{1/2}) \leq D_2(g(n)n)^{1/2} / \varepsilon N.$$

But $g(n) \leq Ne^{-cn/N}$ and $n/N \rightarrow \infty$, so the right-hand side of (3.24) converges to 0.

This proves (3.18) and completes the proof of the theorem.

4. The waiting time for a small fraction. As above let T_b denote the number of balls thrown until exactly $b = b_N$ cells remain empty. Let t_b be the unique solution of the equation

$$(4.1) \quad b = g(t_b) = \sum_{k=1}^N \exp(-t_b p_k).$$

THEOREM 4. *If*

$$(4.2) \quad b_N \rightarrow +\infty,$$

$$(4.3) \quad b_N/N \rightarrow 0,$$

as $N \rightarrow \infty$ and

$$(4.4) \quad 0 < C \leq Np_k \leq D < \infty, \quad \text{for all } k \text{ and } N,$$

then

$$(4.5) \quad b_N^{-1}(T_b - t_b) \sum_{k=1}^N p_k \exp(-t_b p_k) \Rightarrow N(0, 1).$$

PROOF. From (4.1) and (4.4) it follows that

$$(4.6) \quad Cb/N \leq \Delta = \sum_{k=1}^N p_k \exp(-t_b p_k) \leq Db/N.$$

Thus for N sufficiently large

$$(4.7) \quad 0 < C \leq \Delta \cdot N/b \leq D < \infty.$$

As in the proof of Theorem 2 the relation

$$(4.8) \quad P((T_b - t_b)\Delta/b^{1/2} \leq x) = P(Y_n \leq b),$$

holds, where

$$(4.9) \quad n = [t_b + xb^{1/2}/\Delta].$$

We have

$$(4.10) \quad \begin{aligned} g(n)(1 + o(1)) &= g(t_b + xb^{1/2}/\Delta) \\ &= \sum \exp(-t_b p_k) \cdot (1 - xp_k b^{1/2}/\Delta + O(1/b)) \\ &= b - x \cdot b^{1/2} + O(1), \end{aligned}$$

and thus

$$(4.11) \quad g(n) \rightarrow +\infty,$$

and from (3.9)

$$(4.12) \quad f(n) \rightarrow +\infty.$$

Furthermore,

$$(4.13) \quad b = g(t_b) \geq N \exp(-Dt_b/N),$$

so by (4.3) we have

$$(4.14) \quad t_b/N \rightarrow +\infty,$$

and therefore

$$(4.15) \quad n/N \rightarrow +\infty.$$

Hence the assumptions of Theorem 3 are satisfied. Now (4.8) and (4.10) give

$$(4.16) \quad \begin{aligned} P(T_b - t_b) \Delta / b^{\frac{1}{2}} \leq x) \\ &= P(Y_n \leq b) = \Phi((b - g(n))/(g(n))^{\frac{1}{2}}) + o(1) \\ &= \Phi((xb^{\frac{1}{2}} + O(1))/(b(1 + o(1)))^{\frac{1}{2}}) + o(1) \rightarrow \Phi(x), \end{aligned}$$

where $\Phi(x)$ is the standardized normal distribution function. This proves the theorem.

Acknowledgment. I wish to thank Professor B. Harris for pointing out some errors in the original draft of this manuscript and the referee for giving an improved proof of Theorem 3.

REFERENCES

- [1] BAUM, L. E. and BILLINGSLEY, P. (1965). Asymptotic distributions for the coupon collector's problem. *Ann. Math. Statist.* **36** 1835-1839.
- [2] FELLER, W. (1968). *An Introduction to Probability Theory and its Applications* **1**, 3rd ed. Wiley, New York.
- [3] HOLST, L. (1971). Limit theorems for some occupancy and sequential occupancy problems. *Ann. Math. Statist.* **42** 1671-1680.
- [4] KARLIN, S. (1967). Central limit theorems for certain infinite urn schemes. *J. Math. Mech.* **17** 373-401.
- [5] KOLCHIN, V. F. and CHISTYAKOV, V. P. (1974). Combinatorial problems of probability theory. *Itogi Nauki i Tekhniki. Teoriya Veroyatnostei Matematicheskaya Statistika. Teoreticheskaya Kibernetika* **11** 5-45. (Translation in *J. Soviet Mathematics* **4** 217-243.)
- [6] SAMUEL-CAHN, E. (1974). Asymptotic distributions for occupancy and waiting time problems with positive probability of falling through the cells. *Ann. Probability* **2** 515-521.
- [7] SEVASTYANOV, B. A. (1972). Poisson limit law for a scheme of sums of dependent random variables. *Theor. Probability Appl.* **17** 695-699.

DEPARTMENT OF MATHEMATICS
UPPSALA UNIVERSITY
S-75252 UPPSALA, SWEDEN