



THE NUMBER OF K -TONS IN THE COUPON COLLECTOR PROBLEM

JOHN C. SAUNDERS *, Middle Tennessee State University

Abstract

Consider the coupon collector problem where each box of a brand of cereal contains a coupon and there are n different types of coupons. Suppose that the probability of a box containing a coupon of a specific type is $1/n$, and that we keep buying boxes until we collect at least m coupons of each type. For $k \geq m$ call a certain coupon a k -ton if we see it k times by the time we have seen m copies of all of the coupons. Here we determine the asymptotic distribution of the number of k -tons after we have collected m copies of each coupon for any k in a restricted range, given any fixed m . We also determine the asymptotic joint probability distribution over such values of k , and the total number of coupons collected.

Keywords: Probability distributions; combinatorial probability

2020 Mathematics Subject Classification: Primary 60E05
Secondary 60C05

1. Introduction

Consider the coupon collector problem where each box of a brand of cereal contains a coupon and there are n different types of coupons. Suppose that the probability of a box containing a coupon of a specific type is $1/n$, and that we keep buying boxes until we collect at least one coupon of each type. It is well known (see, for example, [9]) that the expected number of boxes we need to buy is nH_n , where H_n is the n th harmonic number:

$$H_n = 1 + \frac{1}{2} + \frac{1}{3} + \cdots + \frac{1}{n} = \log n + \gamma + \frac{1}{2n} + O\left(\frac{1}{n^2}\right),$$

where γ is the Euler–Mascheroni constant. The expected value of the total number of boxes collected has been extensively studied. For instance, it was proved in [9] that, if you continue to collect boxes until you have at least m coupons of each type, then the expected number of boxes collected is $n \log n + (m - 1)n \log \log n + n \cdot C_m + o(n)$, where C_m is a constant depending on m .

Definition 1.1. The Gumbel distribution with parameters μ and β , which is denoted as $\text{Gumbel}(\mu, \beta)$, is defined as the probability distribution with cumulative density function $\exp[-e^{-(x-\mu)/\beta}]$ for $-\infty < x < \infty$.

The above result was improved in [5], which proved that, if $v_m(n)$ is the total number of boxes collected, then $(v_m(n)/n) - \log n - (m - 1) \log \log n \sim \text{Gumbel}(-\log((m - 1)!), 1)$ as

Received 13 August 2021; revision received 7 August 2022.

* Postal address: Middle Tennessee State University Department of Mathematical Sciences. Email: John.Saunders@mtsu.edu

© The Author(s), 2023. Published by Cambridge University Press on behalf of Applied Probability Trust.

$n \rightarrow \infty$. The probability distribution of the number of boxes that need to be collected to collect $a_n + 1$ coupons, where $0 \leq a_n < n$, was studied in [1]. The expected number of boxes needing to be collected if the types of coupons have different probabilities of being in a box was examined in [2], and [3, 8] looked at the actual probability distribution of the number of boxes collected if the types of coupons have different probabilities. The case of collecting at least m coupons of each type has also been examined [4]. Further variations of the problem, where you have multiple collectors, were studied in [6, 7].

Here we examine the situation where we continue to collect boxes until we have at least m copies of each type of coupon for any fixed $m \in \mathbb{N}$.

Definition 1.2. For all $k \geq m$, call a certain coupon a k -ton if we see it k times by the time we have seen m copies of all of the coupons. Let the number of k -tons be denoted as S_k .

In the case of $m = 1$, [7] determined the expected number of 1-tons, which they called singletons, to be H_n . Singletons were also studied in [10], but in the case of unequal coupon probabilities and where the number of coupons collected in total is proportional to the number of types of coupons; many applications to their study, such as in databases, biological particles, and communication channels, were noted here. We extend this singleton result by determining the asymptotic joint probability distribution over values of k in a restricted range and the total number of coupons collected assuming, as in [7], equal coupon probabilities. We also determine the asymptotic distribution of the number of k -tons after we have collected m copies of each coupon for any such values of k , given any fixed $m \in \mathbb{N}$. While making an analogous extension to [10] is certainly interesting, we leave it for a future project.

Theorem 1.1. Fix $m \in \mathbb{N}$, let $T_m(n)$ be the total number of coupons collected to see at least m copies of each coupon, and let $n \rightarrow \infty$. Let $k = o(\log n)$ if $m = 1$ and $k = o((\log n)/(\log \log n))$ if $m \geq 2$, with $k \geq m$ in either case. Then the joint probability distribution of

$$\left(\frac{T_m(n) - n \log n - (m - 1)n \log \log n}{n}, \frac{k!S_k}{(\log n)^{k-m+1}} \right)$$

converges to (X, e^{-X}) , where $X \sim \text{Gumbel}(-\log((m - 1)!), 1)$.

A simple heuristic argument for Theorem 1.1 may be given as follows. First, from [5], we know that

$$\frac{T_m(n) - n \log n - (m - 1)n \log \log n}{n} \sim \text{Gumbel}(-\log((m - 1)!), 1).$$

Later on, we also see, unsurprisingly, that the number of times a specific type of coupon is seen is essentially independent of whether or not all types of coupons have been collected at least m times. Also, unsurprisingly, it turns out that the number of times two different types of coupons are seen is asymptotically independent. Since the probability of a type of coupon being seen k times is

$$\binom{T_m(n)}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{T_m(n)-k} \sim \frac{(n \log n)^k e^{-\log n - (m-1) \log \log n - x}}{k!n^k} = \frac{(\log n)^{k-m+1}}{ne^x k!},$$

where $x = (T_m(n) - n \log n - (m - 1)n \log \log n)/n$, the expected number of k -tons is $\sim (\log n)^{k-m+1}/e^x k!$, so that $k!S_k/(\log n)^{k-m+1} \sim e^{-x}$.

While Theorem 1.1 provides us with the joint distribution between the total number of coupons collected and the number of k -tons for minimal k , we can also ask what happens for maximum values of k . The answer to this is provided in the next theorem.

Theorem 1.2. Fix $m \in \mathbb{N}$ and let $T_m(n)$ be the total number of coupons collected to see at least m copies of each coupon.

- (i) Fix $d \in \mathbb{R}$. Pick an increasing sequence $(n_j) \subset \mathbb{N}$ and a sequence $(d_j) \subset \mathbb{R}$ such that $\lim_{j \rightarrow \infty} d_j = d$ and $k_j = e \log n_j + ((e - 1)(m - 1) - \frac{1}{2}) \log \log n_j + d_j \in \mathbb{N}$ for all $j \in \mathbb{N}$. Then the joint probability distribution of

$$\left(\frac{T_m(n_j) - n_j \log n_j - (m - 1)n_j \log \log n_j}{n_j}, S_{k_j} \right)$$

converges to

$$\left(X, \text{Pois} \left(\frac{e^{(e-1)X-d}}{\sqrt{2\pi e}} \right) \right),$$

where $X \sim \text{Gumbel}(-\log((m - 1)!), 1)$.

- (ii) Let $g(n) = o(\log \log n)$ with $\lim_{n \rightarrow \infty} g(n) = \infty$, such that $k = e \log n + ((e - 1)(m - 1) - \frac{1}{2}) \log \log n + g(n) \in \mathbb{N}$ for all $n \in \mathbb{N}$. Then S_k converges to 0 with probability 1.
- (iii) Let $g(n) = o(\log \log n)$ with $\lim_{n \rightarrow \infty} g(n) = \infty$, such that $k = e \log n + ((e - 1)(m - 1) - \frac{1}{2}) \log \log n - g(n) \in \mathbb{N}$ for all $n \in \mathbb{N}$. Then the value of $\mathbb{P}(S_k = 0)$ converges to 0.

A simple heuristic argument for Theorem 1.2 may be given as follows. The probability of a type of coupon being seen k_j times is

$$\binom{T_m(n_j)}{k_j} \left(\frac{1}{n_j} \right)^{k_j} \left(1 - \frac{1}{n_j} \right)^{T_m(n_j) - k_j}.$$

Assuming that $x = (T_m(n_j) - n_j \log n_j - (m - 1)n_j \log \log n_j) / n_j$ is constant as $j \rightarrow \infty$ and k_j is as defined in Theorem 1.2, we have that the above is asymptotic to

$$\begin{aligned} & \frac{1}{\sqrt{2\pi k_j}} \left(\frac{e T_m(n_j)}{n_j k_j} \right)^{k_j} \exp \left[-\frac{T_m(n_j)}{n_j} \right] \\ & \sim \frac{1}{\sqrt{2\pi e \log n_j}} \exp \left[\frac{T_m(n_j)e - k_j n_j - T_m(n_j)}{n_j} \right] \\ & = \frac{1}{\sqrt{2\pi e \log n_j}} \exp \left[-\log n_j + \frac{\log \log n_j}{2} + (e - 1)x - d_j \right] \\ & = \frac{1}{\sqrt{2\pi e n_j}} e^{(e-1)x - d_j}. \end{aligned}$$

Therefore, by similar reasoning to the heuristic for Theorem 1.1, we can see that the probability that there are y k -tons is asymptotic to

$$\binom{n}{y} \left(\frac{\lambda}{n} \right)^y \left(1 - \frac{\lambda}{n} \right)^{n-y} \sim \frac{\lambda^y e^{-\lambda}}{y!},$$

where $\lambda = e^{(e-1)x-d}/(\sqrt{2\pi e})$, which gives the Poisson distribution in Theorem 1.2(i). Parts (ii) and (iii) can be heuristically argued by letting d_j approach ∞ or $-\infty$.

From Theorem 1.1 we can derive the following corollaries.

Corollary 1.1. Fix $m \in \mathbb{N}$, let $n \rightarrow \infty$, and let $k = o(\log n)$ if $m = 1$ and $k = o((\log n)/(\log \log n))$ if $m \geq 2$, with $k \geq m$ in either case. Suppose we keep collecting coupons until we see at least m copies of each coupon. Then

$$\frac{k!S_k}{(\log n)^{k-m+1}} \xrightarrow{D} \text{Exp}\left(\frac{1}{(m-1)!}\right).$$

Corollary 1.1 gives the asymptotic probability distribution of the number of k -tons. If we fix k , we can also determine the joint asymptotic distribution between the number of m -tons, $m+1$ -tons, \dots , and k -tons. We can see that asymptotically all of these variables are linearly dependent.

Corollary 1.2. Fix $k, m \in \mathbb{N}$ with $k \geq m$, and let $n \rightarrow \infty$. Suppose we keep collecting coupons until we see at least m copies of each coupon. Then the joint probability distribution

$$\left(\frac{m!S_m}{(\log n)}, \frac{(m+1)!S_{m+1}}{(\log n)^2}, \dots, \frac{k!S_k}{(\log n)^{k-m+1}}\right)$$

converges to $X \cdot (1, 1, \dots, 1)$, where

$$X \sim \text{Exp}\left(\frac{1}{(m-1)!}\right).$$

2. Proofs

To prove Theorem 1.1 we require the following notation and lemma.

Note 2.1. Suppose we collect x coupons in total, regardless of whether we have collected all n types of coupons or not. Let $S_{i,x}$ denote the number of types of coupons we collected exactly i times.

Note 2.2. Throughout this paper, let $N(n, f(n)) := n \log n + (m-1)n \log \log n + f(n)$.

Lemma 2.1. Let $f: \mathbb{N} \rightarrow \mathbb{R}$ be such that $-\frac{1}{2}n \log \log n \leq f(n) \leq \frac{1}{2}n \log \log n$, and $N(n, f(n)) \in \mathbb{N}$ for all $n \in \mathbb{N}$. Suppose we collect $N(n, f(n))$ coupons in total, regardless of whether we collect m copies of all n different kinds of coupons or not. Let $k = o(\log n)$ if $m = 1$ and $k = o(\log n / \log \log n)$ if $m \geq 2$. If $m \leq k$, then as $n \rightarrow \infty$

$$\frac{k!S_{k, N(n, f(n))}}{(\log n)^{k-m+1}}$$

is asymptotic to $e^{-f(n)/n}$ with probability 1, and

$$\text{Var}\left(\frac{k!S_{k, N(n, f(n))}}{(\log n)^{k-m+1}}\right) = O\left(\frac{1}{\sqrt{\log n}}\right).$$

Proof. First, assume that $m \leq k$. We can see that, as $n \rightarrow \infty$,

$$\begin{aligned} \mathbb{E}(S_{k,N(n,f(n))}) &= n \cdot \mathbb{P}(c_0 \text{ occurs exactly } k \text{ times in the collection of } N(n, f(n)) \text{ coupons}) \\ &= \binom{N(n, f(n))}{k} \left(1 - \frac{1}{n}\right)^{N(n, f(n))-k} \frac{1}{n^{k-1}} \\ &\sim \frac{N(n, f(n))^k}{k!n^{k-1}} \exp \left[-\log n - (m-1) \log \log n - \frac{f(n)}{n} \right] \\ &\sim \frac{(\log n)^{k-m+1}}{k!e^{f(n)/n}}. \end{aligned}$$

Thus, we can deduce that

$$\mathbb{E} \left(\frac{k!S_{k,N(n,f(n))}}{(\log n)^{k-m+1}} \right) = e^{-f(n)/n} (1 + o(1))$$

as $n \rightarrow \infty$. Also,

$$\begin{aligned} \mathbb{E}(S_{k,N(n,f(n))}^2) &= n(n-1) \cdot \mathbb{P}(c_0, c_1 \text{ both occur exactly } k \text{ times in the collection of } N(n, f(n)) \text{ coupons}) \\ &\quad + n \cdot \mathbb{P}(c_0 \text{ occurs exactly } k \text{ times in the collection of } N(n, f(n)) \text{ coupons}) \\ &= n(n-1) \binom{N(n, f(n))}{k} \binom{N(n, f(n)) - k}{k} \left(1 - \frac{2}{n}\right)^{N(n, f(n))-2k} \frac{1}{n^{2k}} \\ &\quad + \binom{N(n, f(n))}{k} \left(1 - \frac{1}{n}\right)^{N(n, f(n))-k} \frac{1}{n^{k-1}} \\ &\leq (\mathbb{E}(S_{k,N(n,f(n))}))^2 \left(1 - \frac{1}{n}\right)^{-2k} + \mathbb{E}(S_{k,N(n,f(n))}). \end{aligned} \tag{2.1}$$

Notice that

$$\frac{k!e^{-f(n)/n}}{(\log n)^{k-m+1}} \leq \frac{k!}{(\log n)^{k-m+\frac{1}{2}}}. \tag{2.2}$$

If k is bounded as $n \rightarrow \infty$, then the right-hand side of (2.2) is $O(1/\sqrt{\log n})$. On the other hand, for k sufficiently large we have $k! < k^{k-m}$, so even if values of k tend toward ∞ , the right-hand side of (2.2) is $O(1/\sqrt{\log n})$. Thus, from (2.1), we obtain

$$\begin{aligned} &\mathbb{E} \left(\left(\frac{k!S_{k,N(n,f(n))}}{(\log n)^{k-m+1}} \right)^2 \right) - \left(\mathbb{E} \left(\frac{k!S_{k,N(n,f(n))}}{(\log n)^{k-m+1}} \right) \right)^2 \\ &\leq \left(\mathbb{E} \left(\frac{k!S_{k,N(n,f(n))}}{(\log n)^{k-m+1}} \right) \right)^2 \cdot O \left(\frac{k}{n} \right) + \frac{(k!)^2 \mathbb{E}(S_{k,N(n,f(n))})}{(\log n)^{2k-2m+2}} \\ &\leq O \left(\frac{k\sqrt{\log n}}{n} \right) + O \left(\frac{1}{\sqrt{\log n}} \right) = O \left(\frac{1}{\sqrt{\log n}} \right) \end{aligned}$$

as $n \rightarrow \infty$. Thus,

$$\text{Var}\left(\frac{k!S_{k,N(n,f(n))}}{(\log n)^{k-m+1}}\right) = O\left(\frac{1}{\sqrt{\log n}}\right)$$

as $n \rightarrow \infty$, so that we obtain our result for $m \leq k$. \square

Proof of Theorem 1.1. $X = (T_m(n) - n \log n - (m-1)n \log \log n)/n$ was proved to be asymptotically Gumbel distributed in [5]. To prove Theorem 1.1, it therefore suffices to show that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{T_m(n) - n \log n - (m-1)n \log \log n}{n} < x, \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) = 0 \quad (2.3)$$

for any fixed $x < y$, and that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(\frac{T_m(n) - n \log n - (m-1)n \log \log n}{n} > x, \frac{k!S_k}{(\log n)^{k-m+1}} > e^{-y}\right) = 0 \quad (2.4)$$

for any fixed $y < x$. To prove (2.3), we first calculate an upper bound for

$$\mathbb{P}\left(T_m(n) \leq N(n, nx), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right)$$

for all sufficiently large $n \in \mathbb{N}$. Let $M(n) = n \log n + (m - \frac{3}{2})n \log \log n$. We have

$$\begin{aligned} & \mathbb{P}\left(T_m(n) \leq N(n, nx), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &= \mathbb{P}\left(T_m(n) < M(n), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &+ \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &\leq \mathbb{P}(T_m(n) < M(n)) + \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right). \end{aligned}$$

We have

$$\begin{aligned} & \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &= \sum_{M=\lceil M(n) \rceil}^{\lfloor N(n, nx) \rfloor} \mathbb{P}\left(T_m(n) = M, \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right). \quad (2.5) \end{aligned}$$

Clearly,

$$\begin{aligned} & \mathbb{P}\left(T_m(n) = M, \frac{k!S_k}{(\log n)^{k-m+1}} < e^{-y}\right) \leq \mathbb{P}\left(S_{m-1, M-1} = 1, S_{m-1, M} = 0, \frac{k!S_{k, M}}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &= \frac{1}{n} \mathbb{P}\left(S_{m-1, M-1} = 1, \frac{k!S_{k, M}}{(\log n)^{k-m+1}} < e^{-y}\right) \\ &\leq \frac{1}{n} \mathbb{P}\left(\frac{k!S_{k, M}}{(\log n)^{k-m+1}} < e^{-y}\right). \quad (2.6) \end{aligned}$$

Combining (2.5) and (2.6), we therefore have

$$\begin{aligned} \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_{k,M}}{(\log n)^{k-m+1}} < e^{-y}\right) \\ \leq \sum_{M=\lceil M(n) \rceil}^{\lfloor N(n, nx) \rfloor} \frac{1}{n} \mathbb{P}\left(\frac{k!S_{k,M}}{(\log n)^{k-m+1}} < e^{-y}\right). \end{aligned} \tag{2.7}$$

From Chebyshev’s inequality and Lemma 2.1 we have, for sufficiently large n ,

$$\mathbb{P}\left(\frac{k!S_{k,M}}{(\log n)^{k-m+1}} < e^{-y}\right) \leq \frac{C}{\sqrt{\log n}}$$

for some constant $C > 0$, depending only on x and y . From (2.7), we can thus deduce that

$$\begin{aligned} \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_{k,M}}{(\log n)^{k-m+1}} < e^{-y}\right) \\ \leq \sum_{M=\lceil M(n) \rceil}^{\lfloor N(n, nx) \rfloor} \frac{C}{n\sqrt{\log n}} \\ \leq \frac{C(nx + \frac{1}{2}n \log \log n)}{n\sqrt{\log n}} = \frac{C(2x + \log \log n)}{2\sqrt{\log n}} \end{aligned} \tag{2.8}$$

for all sufficiently large n .

Now we let $n \rightarrow \infty$. From the result quoted at the start of the proof, we can deduce that $\lim_{n \rightarrow \infty} \mathbb{P}(T < M(n)) = 0$. From (2.8) we can also see that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left(M(n) \leq T_m(n) \leq N(n, nx), \frac{k!S_{k,M}}{(\log n)^{k-m+1}} < e^{-y}\right) = 0.$$

These two limits give us (2.3).

We can prove (2.4) similarly, replacing $M(n) = n \log n + (m - \frac{3}{2})n \log \log n$ with $M'(n) = n \log n + (m - \frac{1}{2})n \log \log n$ and reversing the appropriate inequalities. \square

Proof of Corollaries 1.1 and 1.2. Corollaries 1.1 and 1.2 follow from Theorem 1.1 by showing that, if $X \sim \text{Gumbel}(-\log((m-1)!), 1)$, then

$$e^{-X} \sim \text{Exp}\left(\frac{1}{(m-1)!}\right).$$

Indeed, we have, for any $r > 0$,

$$\begin{aligned} \mathbb{P}(e^{-X} \leq r) &= \mathbb{P}(X \geq -\log r) = 1 - \mathbb{P}(X \leq -\log r) \\ &= 1 - \exp\left[-e^{-(-\log r + \log(m-1)!)}\right] \\ &= 1 - e^{-r/(m-1)!}. \end{aligned}$$

Hence, the result follows. \square

To prove Theorem 1.2, we require the following lemma.

Lemma 2.2. Fix $m \in \mathbb{N}$. Let $f: \mathbb{N} \rightarrow \mathbb{R}$ be such that there exists $x \in \mathbb{R}$ such that $\lim_{n \rightarrow \infty} f(n)/n = x$ and $N(n, f(n)) \in \mathbb{N}$ for all $n \in \mathbb{N}$. Suppose we collect $N(n, f(n))$ coupons in total regardless of whether we collect m copies of all n different kinds of coupons or not. Pick an increasing sequence $(n_j)_j \subset \mathbb{N}$ and a sequence $(d_j)_j \subset \mathbb{R}$ such that $\lim_{j \rightarrow \infty} d_j = d$ for some $d \in \mathbb{R}$ and $k = e \log n_j + ((e-1)(m-1) - \frac{1}{2}) \log \log n_j + d_j \in \mathbb{N}$ for all $j \in \mathbb{N}$. Then, as $j \rightarrow \infty$,

$$\begin{aligned} \mathbb{P}(S_{m-1, N(n_j, f(n_j))} = r_1, S_{k, N(n_j, f(n_j))} = r_2) &= \frac{1}{r_1! r_2!} \left(\frac{e^{-f(n_j)/n_j}}{(m-1)!} \right)^{r_1} \exp \left[\frac{-e^{-f(n_j)/n_j}}{(m-1)!} \right] \\ &\times \left(\frac{\exp[(e-1)f(n_j)/n_j - d]}{\sqrt{2\pi e}} \right)^{r_2} \\ &\times \exp \left[\frac{-\exp[(e-1)f(n_j)/n_j - d]}{\sqrt{2\pi e}} \right] \\ &\times (1 + o(1)), \end{aligned}$$

where the implied constant in the error term only depends upon r_1 and r_2 .

Proof. Take some $j \in \mathbb{N}$. Let the probability that r_1 prescribed types of coupons are $(m-1)$ -tons and r_2 prescribed types of coupons are k -tons be denoted as $W_{r_1, r_2}(n)$. Then, as $j \rightarrow \infty$,

$$\begin{aligned} &\binom{n_j}{r_1, r_2, n_j - r_1 - r_2} W_{r_1, r_2}(n_j) \\ &= \binom{n_j}{r_1, r_2, n - r_1 - r_2} \frac{N(n_j, f(n_j))!}{(m-1)!^{r_1} k!^{r_2} (N(n_j, f(n_j)) - (m-1)r_1 - kr_2)! n_j^{(m-1)r_1 + kr_2}} \\ &\quad \times \left(1 - \frac{r_1 + r_2}{n_j} \right)^{N(n_j, f(n_j)) - (m-1)r_1 - kr_2} \\ &= \binom{n_j}{r_1, r_2, n_j - r_1 - r_2} \frac{N(n_j, f(n_j))! e^{-(r_1 + r_2)N(n_j, f(n_j))/n_j}}{(m-1)!^{r_1} k!^{r_2} (N(n_j, f(n_j)) - (m-1)r_1 - kr_2)! n_j^{(m-1)r_1 + kr_2}} \\ &\quad \times \left(1 + O\left(\frac{\log n_j}{n_j} \right) \right) \\ &= \frac{N(n_j, f(n_j))^{(m-1)r_1 + kr_2} e^{-(r_1 + r_2)f(n_j)/n_j}}{r_1! r_2! (m-1)!^{r_1} k!^{r_2} (\log n_j)^{(r_1 + r_2)(m-1)} n_j^{(m-1)r_1 + kr_2}} \left(1 + O\left(\frac{\log n_j}{n_j} \right) \right) \\ &= \frac{1}{r_1! r_2!} \left(\frac{e^{-f(n_j)/n_j}}{(m-1)!} \right)^{r_1} \left(\frac{e^{((e-1)f(n_j)/n_j) - d_j}}{\sqrt{2\pi e}} \right)^{r_2} \left(1 + O\left(\frac{(\log \log n_j)^2}{\log n_j} \right) \right), \end{aligned}$$

where the implied constant in the error term only depends on r_1 and r_2 . By the inclusion–exclusion formula, we have

$$\begin{aligned} & \mathbb{P}(S_{m-1, N(n_j, f(n_j))} = r_1, S_{k, N(n_j, f(n_j))} = r_2) \\ &= \sum_{0 \leq j_1 + j_2 \leq n_j - r_1 - r_2} \left[(-1)^{j_1 + j_2} \binom{n_j}{r_1 + j_1, r_2 + j_2, n_j - r_1 - r_2 - j_1 - j_2} \binom{r_1 + j_1}{r_1} \binom{r_2 + j_2}{r_2} \right. \\ & \qquad \qquad \qquad \left. \times W_{r_1 + j_1, r_2 + j_2}(n_j) \right]. \end{aligned} \tag{2.9}$$

For $t \leq n_j - r_1 - r_2$ even:

$$\begin{aligned} & \mathbb{P}(S_{m-1, N(n_j, f(n_j))} = r_1, S_{k, N(n_j, f(n_j))} = r_2) \\ & \leq \sum_{0 \leq j_1 + j_2 \leq t} \left[(-1)^{j_1 + j_2} \binom{n_j}{r_1 + j_1, r_2 + j_2, n_j - r_1 - r_2 - j_1 - j_2} \binom{r_1 + j_1}{r_1} \binom{r_2 + j_2}{r_2} \right. \\ & \qquad \qquad \qquad \left. \times W_{r_1 + j_1, r_2 + j_2}(n_j) \right]. \end{aligned} \tag{2.10}$$

For $t \leq n_j - r_1 - r_2$ odd:

$$\begin{aligned} & \mathbb{P}(S_{m-1, N(n_j, f(n_j))} = r_1, S_{k, N(n_j, f(n_j))} = r_2) \\ & \geq \sum_{0 \leq j_1 + j_2 \leq t} \left[(-1)^{j_1 + j_2} \binom{n_j}{r_1 + j_1, r_2 + j_2, n_j - r_1 - r_2 - j_1 - j_2} \binom{r_1 + j_1}{r_1} \binom{r_2 + j_2}{r_2} \right. \\ & \qquad \qquad \qquad \left. \times W_{r_1 + j_1, r_2 + j_2}(n_j) \right]. \end{aligned} \tag{2.11}$$

Combining (2.9)–(2.11), we have

$$\begin{aligned} & \mathbb{P}(S_{m-1, N(n_j, f(n_j))} = r_1, S_{k, N(n_j, f(n_j))} = r_2) \\ & \sim \frac{1}{r_1! r_2!} \sum_{j_1, j_2=0}^{\infty} \frac{(-1)^{j_1 + j_2}}{j_1! j_2!} \left(\frac{e^{-x}}{(m-1)!} \right)^{r_1 + j_1} \left(\frac{e^{((e-1)x-d)}}{\sqrt{2\pi e}} \right)^{r_2 + j_2} \\ & = \frac{1}{r_1! r_2!} \left(\frac{e^{-x}}{(m-1)!} \right)^{r_1} \exp \left[\frac{-e^{-x}}{(m-1)!} \right] \left(\frac{e^{((e-1)x-d)}}{\sqrt{2\pi e}} \right)^{r_2} \exp \left[\frac{-e^{((e-1)x-d)}}{\sqrt{2\pi e}} \right], \end{aligned}$$

as $j \rightarrow \infty$. □

Proof of Theorem 1.2. We first prove Theorem 1.2(i). Let $d \in \mathbb{R}$ and pick increasing sequences $(n_j)_j$ and $(d_j)_j$ such that $\lim_{j \rightarrow \infty} d_j = d$ and $k_j := e \log n_j + ((e-1)(m-1) - \frac{1}{2}) \log \log n_j + d_j \in \mathbb{N}$ for all $j \in \mathbb{N}$. Let $x < y$ and $r \in \mathbb{N}$ be constants. Consider the probability $\mathbb{P}(N(n_j, n_j x) \leq T_m(n) \leq N(n_j, n_j y), S_k = r)$. We have

$$\mathbb{P}(N(n_j, n_j x) \leq T_m(n) \leq N(n_j, n_j y), S_k = r) = \sum_{M=\lceil N(n_j, n_j x) \rceil}^{\lfloor N(n_j, n_j y) \rfloor} \mathbb{P}(T = M, S_{k, M-1} = r).$$

Clearly,

$$\begin{aligned} \mathbb{P}(T_m(n) = M, S_{k,M-1} = r) &\leq \mathbb{P}(S_{m-1,M-1} = 1, S_{m-1,M} = 0, S_{k,M-1} = r) \\ &= \mathbb{P}(S_{m-1,M-1} = 1, S_{k,M-1} = r)/n. \end{aligned} \tag{2.12}$$

Also, for $M \geq N(n_j, n_jx)$, if we have $S_{i,M} > 0$ for some $0 \leq i \leq m - 2$, then we have $S_{l, \lceil N(n_j, n_jx) \rceil} > 0$ for some $0 \leq l \leq i$. Thus, we can also deduce that

$$\begin{aligned} \mathbb{P}(N(n_j, n_jx) \leq T_m(n) \leq N(n_j, n_jy), S_k = r) \\ \geq \left(\sum_{M=\lceil N(n_j, n_jx) \rceil}^{\lfloor N(n_j, n_jy) \rfloor} \frac{\mathbb{P}(S_{m-1,M-1} = 1, S_{k,M-1} = r)}{n_j} \right) \\ - \mathbb{P}(\text{there exists } 0 \leq l \leq m - 2 \text{ such that } S_{l, \lceil N(n_j, n_jx) \rceil} > 0). \end{aligned} \tag{2.13}$$

For any $0 \leq l \leq m - 2$, notice that, as $j \rightarrow \infty$,

$$\begin{aligned} \mathbb{E}(S_{l, \lceil N(n_j, n_jx) \rceil}) &= n \cdot \mathbb{P}(c_0 \text{ occurs exactly } l \text{ times in the collection of } \lceil N(n_j, n_jx) \rceil \text{ coupons}) \\ &= \binom{\lceil N(n_j, n_jx) \rceil}{l} \left(1 - \frac{1}{n_j}\right)^{\lceil N(n_j, n_jx) \rceil - l} \frac{1}{n_j^{l-1}} \\ &\sim \frac{N(n_j, n_jx)^l}{l! n_j^{l-1}} e^{-\log n_j - (m-1) \log \log n_j - x} \sim \frac{(\log n_j)^{l-m+1} e^{-x}}{l!}. \end{aligned}$$

Since $l \leq m - 2$, we therefore have $\lim_{j \rightarrow \infty} \mathbb{E}(S_{l, \lceil N(n_j, n_jx) \rceil}) = 0$, and so

$$\lim_{j \rightarrow \infty} \mathbb{P}(\text{there exists } 0 \leq l \leq m - 2 \text{ such that } S_{l, \lceil N(n_j, n_jx) \rceil} > 0) = 0. \tag{2.14}$$

For every $N(n_j, n_jx) \leq M \leq N(n_j, n_jy)$, let $c_{M,j} := (M - n_j \log n_j - (m - 1)n_j \log \log n_j)/n_j$. By Lemma 2.2, we have

$$\begin{aligned} &\sum_{M=\lceil N(n_j, n_jx) \rceil}^{\lfloor N(n_j, n_jy) \rfloor} \frac{\mathbb{P}(S_{m-1,M-1} = 1, S_{k,M-1} = r)}{n} \\ &= \sum_{M=\lceil N(n_j, n_jx) \rceil}^{\lfloor N(n_j, n_jy) \rfloor} \frac{1}{r! n} \left(\frac{e^{-c_{M,j}}}{(m-1)!} \right) \exp \left[\frac{-e^{-c_{M,j}}}{(m-1)!} \right] \left(\frac{e^{(e-1)c_{M,j}-d}}{\sqrt{2\pi e}} \right)^r \exp \left[\frac{-e^{(e-1)c_{M,j}-d}}{\sqrt{2\pi e}} \right] \\ &\quad \times (1 + o(1)) \\ &< \left(\frac{e^{(e-1)y-d}}{\sqrt{2\pi e}} \right)^r \exp \left[\frac{-e^{(e-1)x-d}}{\sqrt{2\pi e}} \right] \\ &\quad \times \sum_{M=\lceil N(n_j, n_jx) \rceil}^{\lfloor N(n_j, n_jy) \rfloor} \frac{1}{r! n} \left(\frac{e^{-c_{M,j}}}{(m-1)!} \right) \exp \left[\frac{-e^{-c_{M,j}}}{(m-1)!} \right] (1 + o(1)). \end{aligned} \tag{2.15}$$

Note that

$$\begin{aligned}
 \lim_{j \rightarrow \infty} \sum_{M=\lceil N(n_j, n_j x) \rceil}^{\lfloor N(n_j, n_j y) \rfloor} \frac{1}{n} \left(\frac{e^{-cM_j}}{(m-1)!} \right) \exp \left[\frac{-e^{-cM_j}}{(m-1)!} \right] \\
 = \int_x^y \frac{e^{-t} \exp \left[\frac{-e^{-t}}{(m-1)!} \right]}{(m-1)!} dt \\
 = \exp \left[\frac{-e^{-y}}{(m-1)!} \right] - \exp \left[\frac{-e^{-x}}{(m-1)!} \right] \\
 = \lim_{j \rightarrow \infty} \mathbb{P}(N(n_j, x) \leq T_m(n_j) \leq N(n_j, y)). \tag{2.16}
 \end{aligned}$$

Combining (2.12)–(2.16), we obtain

$$\begin{aligned}
 \limsup_{j \rightarrow \infty} \mathbb{P}(N(n_j, n_j x) \leq T_m(n_j) \leq N(n_j, n_j y), S_k = r) \\
 \leq \frac{1}{r!} \left(\frac{e^{(e-1)y-d}}{\sqrt{2\pi e}} \right)^r \exp \left[\frac{-e^{(e-1)x-d}}{\sqrt{2\pi e}} \right] \lim_{j \rightarrow \infty} \mathbb{P}(N(n_j, n_j x) \leq T_m(n_j) \leq N(n_j, n_j y)).
 \end{aligned}$$

By similar reasoning, we can also obtain

$$\begin{aligned}
 \liminf_{j \rightarrow \infty} \mathbb{P}(N(n_j, n_j x) \leq T_m(n_j) \leq N(n_j, n_j y), S_k = r) \\
 \geq \frac{1}{r!} \left(\frac{e^{(e-1)x-d}}{\sqrt{2\pi e}} \right)^r \exp \left[\frac{-e^{(e-1)y-d}}{\sqrt{2\pi e}} \right] \lim_{j \rightarrow \infty} \mathbb{P}(N(n_j, n_j x) \leq T_m(n_j) \leq N(n_j, n_j y)).
 \end{aligned}$$

Since the choices of $x < y$ were arbitrary, we can see that the joint limiting distribution holds. □

For Theorem 1.2(ii), again let $x < y$ be constants. We will show that

$$\lim_{n \rightarrow \infty} \mathbb{P}(N(n, nx) \leq T_m(n) \leq N(n, ny), S_k \geq 1) = 0,$$

which implies the desired result since $x < y$ are arbitrary. We have

$$\begin{aligned}
 & \mathbb{P}(N(n, nx) \leq T_m(n) \leq N(n, ny), S_k \geq 1) \\
 & \leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{\mathbb{P}(S_{m-1, M-1} = 1, S_{k, M-1} \geq 1)}{n} \\
 & \leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{\mathbb{P}(S_{k, M-1} \geq 1)}{n} \\
 & \leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{n \cdot \mathbb{P}(c_0 \text{ occurs exactly } k \text{ times in the collection of } M-1 \text{ coupons})}{n}
 \end{aligned}$$

$$\begin{aligned}
&= \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \mathbb{P}(c_0 \text{ occurs exactly } k \text{ times in the collection of } M \text{ coupons}) \\
&= \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \binom{M}{k} \left(1 - \frac{1}{n}\right)^{M-k} \frac{1}{n^k} \\
&< \frac{1}{k!} \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \left(1 - \frac{1}{n}\right)^{M-k} \frac{M^k}{n^k} \\
&< \frac{n(y-x)}{k!} \left(1 - \frac{1}{n}\right)^{N(n, nx)-k} (\log n + (m-1) \log \log n + y)^k.
\end{aligned}$$

As $n \rightarrow \infty$, we have

$$\begin{aligned}
&\frac{n(y-x)}{k!} \left(1 - \frac{1}{n}\right)^{N(n, nx)-k} (\log n + (m-1) \log \log n + y)^k \\
&\sim \frac{(y-x)(\log n + (m-1) \log \log n + y)^k}{k!(\log n)^{m-1} e^x} \\
&\sim \frac{y-x}{\sqrt{2\pi e}(\log n)^{m-\frac{1}{2}} e^x} \left(\frac{e \log n + e(m-1) \log \log n + ey}{k}\right)^k \\
&\sim \frac{(y-x)e^{(m-\frac{1}{2}) \log \log n + ey - g(n)}}{\sqrt{2\pi e}(\log n)^{m-\frac{1}{2}} e^x} \\
&= \frac{(y-x)e^{ey-x-g(n)}}{\sqrt{2\pi e}}. \tag{2.17}
\end{aligned}$$

We have

$$\lim_{n \rightarrow \infty} \frac{(y-x)e^{ey-x-g(n)}}{\sqrt{2\pi e}} = 0,$$

giving us our result.

For Theorem 1.2(iii), again let $x < y$ be constants. We will show that

$$\lim_{n \rightarrow \infty} \mathbb{P}(N(n, nx) \leq T_m(n) \leq N(n, ny), S_k = 0) = 0,$$

which implies the desired result since $x < y$ are arbitrary. We have

$$\begin{aligned}
\mathbb{P}(N(n, nx) \leq T_m(n) \leq N(n, ny), S_k = 0) &\leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{\mathbb{P}(S_{m-1, M-1} = 1, S_{k, M-1} = 0)}{n} \\
&\leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{\mathbb{P}(S_{k, M-1} = 0)}{n}. \tag{2.18}
\end{aligned}$$

Note that

$$\begin{aligned} \mathbb{E}(S_{k,M}) &= \binom{M}{k} \left(1 - \frac{1}{n}\right)^{M-k} \frac{1}{n^{k-1}}, \\ \mathbb{E}(S_{k,M}^2) &= \binom{M}{k} \binom{M-k}{k} \left(1 - \frac{2}{n}\right)^{M-2k} \frac{(n-1)}{n^{2k-1}} + \binom{M}{k} \left(1 - \frac{1}{n}\right)^{M-k} \frac{1}{n^{k-1}} \\ &< \binom{M}{k}^2 \left(1 - \frac{1}{n}\right)^{2M-4k} \frac{1}{n^{2k-2}} + \binom{M}{k} \left(1 - \frac{1}{n}\right)^{M-k} \frac{1}{n^{k-1}}. \end{aligned}$$

We have

$$\text{Var}(S_{k,M}) = \mathbb{E}(S_{k,M}^2) - \mathbb{E}(S_{k,M})^2 \leq \mathbb{E}(S_{k,M})^2 \left(\left(1 - \frac{1}{n}\right)^{-2k} - 1 \right) + \mathbb{E}(S_{k,M}).$$

Since $k < n$ for sufficiently large n , we have that there exists $C_1 > 0$ such that, for sufficiently large n ,

$$\text{Var}(S_{k,M}) \leq \frac{C_1 k \mathbb{E}(S_{k,M})^2}{n} + \mathbb{E}(S_{k,M}).$$

Similarly to the derivation of (2.17), we can deduce that there exists $C_2 > 0$ depending only on x and y such that, for sufficiently large n , $\mathbb{E}(S_{k,M}) < C_2 e^{g(n)} < C_2 \log n$. Therefore, for sufficiently large n , we have $\text{Var}(S_{k,M}) \leq \mathbb{E}(S_{k,M}) + 1 < 2\mathbb{E}(S_{k,M})$. Thus, for sufficiently large n , $\sigma(S_{k,M}) < \sqrt{2\mathbb{E}(S_{k,M})}$. By Chebyshev’s inequality, we therefore have

$$\begin{aligned} \mathbb{P}(S_{k,M} = 0) &\leq \mathbb{P}(|S_{k,M} - \mathbb{E}(S_{k,M})| \geq \mathbb{E}(S_{k,M})) \\ &\leq \mathbb{P}\left(|S_{k,M} - \mathbb{E}(S_{k,M})| \geq \frac{\sigma(S_{k,M})\sqrt{\mathbb{E}(S_{k,M})}}{\sqrt{2}}\right) \leq \frac{2}{\mathbb{E}(S_{k,M})} \end{aligned} \tag{2.19}$$

for sufficiently large n . Also, similarly to the derivation of (2.17), we have that there exists $C_3 > 0$ depending on only x and y such that $C_3 e^{g(n)} < \mathbb{E}(S_{k,M})$. Thus, from (2.18) and (2.19), we have

$$\mathbb{P}(N(n, nx) \leq T_m(n) \leq N(n, ny), S_k = 0) \leq \sum_{M=\lceil N(n, nx) \rceil}^{\lfloor N(n, ny) \rfloor} \frac{2}{C_3 n e^{g(n)}} = \frac{2(y-x)}{C_3 e^{g(n)}}$$

for sufficiently large n . Since $\lim_{n \rightarrow \infty} g(n) = \infty$, we have our result. □

3. Future work

There are a few open questions that are worth exploring with respect to the coupon collector’s problem and the number of k -tons. For instance, can we extend our ranges for k ? We have results for $k = o(\log n)$ if $m = 1$ and $k = o(\log n / \log \log n)$ if $m \geq 2$, as well as $k = e \log n + ((e - 1)(m - 1) - \frac{1}{2}) \log \log n + d$, but we can also ask similar questions for values of k between these minimum and maximum values. For instance, we can consider $k = \theta(\log n)$ or $k = \theta(\log n / \log \log n)$. We can also ask what happens if we have m increasing with n . Also, [3] studied the limiting distribution of $T_m(n)$ in the case of unequal coupon probabilities. We could also study the limiting distribution of k -tons in the case of unequal coupon probabilities.

Acknowledgement

The author would like to thank Dr. Daniel Berend for helpful comments on the paper.

Funding information

The research of J. C. Saunders is supported by an Azrieli International Postdoctoral Fellowship, as well as a Postdoctoral Fellowship at the University of Calgary.

Competing interests

There were no competing interests to declare which arose during the preparation or publication process of this article.

References

- [1] BAUM, L. AND BILLINGSLEY, P. (1965). Asymptotic distributions for the coupon collector's problem. *Ann. Math. Statist.* **36**, 1835–1839.
- [2] BERENBRINK, P. AND SAUERWALD, T. (2009). The weighted coupon collector's problem and applications. In *Computing and Combinatorics* (LNCS 5609), ed H. Q. Ngo. Springer, Berlin, pp. 449–458.
- [3] DOUMAS, A. AND PAPANICOLAOU, V. (2012). The coupon collector's problem revisited: Asymptotics of the variance. *Adv. Appl. Prob.* **44**, 166–195.
- [4] DOUMAS, A. AND PAPANICOLAOU, V. (2016). The coupon collector's problem revisited: Generalizing the double Dixie cup problem of Newman and Shepp. *ESIAM Prob. Statist.* **20**, 367–399.
- [5] ERDŐS, P. AND RÉNYI, A. (1961). On a classical problem of probability theory. *Magyar Tudományos Akadémia Matematikai Kutató Intézetének Közleményei* **6**, 215–220.
- [6] FOATA, D. AND HAN, G. (2001). Les nombres hyperharmoniques et la fratrie du collectionneur de vignettes. *Séminaire Lotharingien de Combinatoire* **47**, B47a.
- [7] MYERS, A. AND WILF, H. (2006). Some new aspects of the coupon collector's problem. *SIAM Rev.* **48**, 549–565.
- [8] NEAL, P. (2008). The generalised coupon collector problem. *J. Appl. Prob.* **45**, 621–629.
- [9] NEWMAN, D. AND SHEPP, L. (1960). The double Dixie cup problem. *Amer. Math. Monthly* **67**, 58–61.
- [10] PENROSE, M. (2009). Normal approximation for isolated balls in an urn allocation model. *Electron. J. Probab.* **14**, 2155–2181.