# Asymptotics of the rising moments for the coupon collector's problem

Aristides V. Doumas*        Vassilis G. Papanicolaou†

### Abstract

We develop techniques of computing the asymptotics of the so–called rising moments of the number $T_N$ of coupons that a collector has to buy in order to find all $N$ existing different coupons as $N \to \infty$. The probabilities (occurring frequencies) of the coupons can be quite arbitrary. After mentioning the case where the coupon probabilities are equal we consider the general case (of unequal probabilities). For a large class of families of coupon probabilities, after adopting a dichotomy, we arrive at the leading behavior of the rising moments of $T_N$ as $N \to \infty$. We also present various illustrative examples.

## 1   Introduction - History of the problem

Consider a population whose members are of $N$ different *types* (e.g. the population may consist of fish, viruses, English words, or baseball cards). For $1 \leq j \leq N$ we denote by $p_j$ the probability that a member of the population is of type $j$. The members of the population are sampled independently with replacement and their types are recorded. The so-called **"Coupon Collector's Problem" (CCP)** deals with questions arising in the above procedure. In particular, CCP pertains to the family of urn problems. Other classical such problems are birthday, Dixie cup or occupancy problems, whose origin can be traced back to De Moivre's treatise *De Mensura Sortis* of 1712 (see, e.g., [16]). CCP (in its simplest form, i.e. the case of equal probabilities) had appeared in W. Feller's classical work [12] and has attracted the attention of various researchers, since it has found many applications in many areas of science (computer science–search algorithms, mathematical programming, optimization, learning processes, engineering, ecology, as well as linguistics, — see, e.g., [5]). Let $T_N$ be the number of trials it takes until all $N$ types are detected (at least once). Apart from its distribution some other interesting quantities are the moments (or, equivalently, the rising moments) of the random variable $T_N$. For the case of equal sampling probabilities the first and the second moment of $T_N$

---

*National Technical University of Athens, Greece. E-mail: `aris.doumas@hotmail.com`
†National Technical University of Athens, Greece. E-mail: `papanico@math.ntua.gr`

are well known and, furthermore, asymptotics as well as limiting results have been obtained by several authors (see for instance [19], [11], [3], [17], [16], [10], [18], and [8]). In particular, in [19] the authors answered the question: how long, on average does it take to obtain $m$ complete sets of $N$ coupons. For unequal probabilities, general asymptotic estimates regarding the first and the second moment, as well as for the variance, have also been obtained by several authors (see, e.g., [7], [13], [8], [9]).

Let $r \geq 1$, be an integer. Set

$$T_N^{(r)} := T_N \left(T_N + 1\right) \left(T_N + 2\right) \cdots \left(T_N + r - 1\right). \tag{1.1}$$

In this paper we consider the $r$-th rising moment of $T_N$, namely

$$E\left[T_N^{(r)}\right] = E\left[T_N \left(T_N + 1\right) \left(T_N + 2\right) \cdots \left(T_N + r - 1\right)\right]. \tag{1.2}$$

In Section 2 we present general expressions for $E\left[T_N^{(r)}\right]$ and exhibit well-known results, mainly for the simplest case of the problem, i.e. the case of equal probabilities. We also describe the general setup of the problem considered in the present paper. In Section 3 we begin by discussing a key feature, namely that the families of the coupon probabilities, i.e. the $p_j$'s, can be divided in two types. The main result for the $p_j$'s of the first type is presented in Theorem 3.5 (the so-called *linear* case falls in this category). Then, we consider a large class of families of coupon probabilities belonging to the second type. The (leading) asymptotic behavior of the rising moments $E\left[T_N^{(r)}\right]$ is given in Theorem 3.7 (the generalized *Zipf* law falls in this case). Furthermore, Theorem 3.9 helps us obtain asymptotic estimates by comparison with cases for which the asymptotic estimates are known. Section 4 contains various examples. Finally, we mention some possible extensions at the end of the paper.

## 2  Preliminaries

For each $j \in \{1, ..., N\}$ it is convenient to introduce the event $A_j^k$, that the type $j$ is **not** detected until trial $k$ (included). Then

$$P\{T_N \geq k\} = P\left(A_1^{k-1} \cup \cdots \cup A_N^{k-1}\right), \qquad k = 1, 2, \dots.$$

By invoking the inclusion-exclusion principle one gets

$$P\{T_N \geq k\} = \sum_{\substack{J \subset \{1,...,N\} \\ J \neq \emptyset}} (-1)^{|J|-1} \left[1 - \left(\sum_{j \in J} p_j\right)\right]^{k-1}, \qquad k = 1, 2, ..., \tag{2.1}$$

where the sum extends over all $2^N - 1$ nonempty subsets $J$ of $\{1, ..., N\}$, while $|J|$ denotes the cardinality of $J$. For $z \in \mathbb{C}, |z| \geq 1$, we introduce the following moment generating function of $T_N$,

$$G(z) := E\left[z^{-T_N}\right] = 1 + \left(z^{-1} - 1\right) \sum_{k=1}^{\infty} z^{-(k-1)} P\{T_N \geq k\} \tag{2.2}$$

(the derivation of the second equality is based on Abel's partial summation formula). Consequently, by using (2.1) one arrives at

$$G(z) = 1 + \left(z^{-1} - 1\right) \sum_{\substack{J \subset \{1,...,N\} \\ J \neq \emptyset}} (-1)^{|J|-1} \sum_{k=1}^{\infty} z^{-(k-1)} \left[1 - \left(\sum_{j \in J} p_j\right)\right]^{k-1}$$

and, hence, by summing the geometric series

$$G(z) = 1 - (z-1) \sum_{\substack{J \subset \{1,\dots,N\} \\ J \neq \emptyset}} \frac{(-1)^{|J|-1}}{z-1+\left(\sum_{j \in J} p_j\right)}. \tag{2.3}$$

We proceed by noticing that

$$\prod_{j=1}^{N} \left(1 - e^{-p_j t}\right) = \sum_{J \subset \{1,\dots,N\}} (-1)^{|J|} \exp\left(-t \sum_{j \in J} p_j\right). \tag{2.4}$$

Thus, at least for $\Re\{z\} \geq 1$,

$$\int_0^\infty \left[1 - \prod_{j=1}^{N} \left(1 - e^{-p_j t}\right)\right] e^{-(z-1)t} dt = \sum_{\substack{J \subset \{1,\dots,N\} \\ J \neq \emptyset}} \frac{(-1)^{|J|-1}}{z-1+\left(\sum_{j \in J} p_j\right)}. \tag{2.5}$$

Finally, by comparing (2.3) and (2.5) we get

$$G(z) = 1 - (z-1) \int_0^\infty \left[1 - \prod_{j=1}^{N} \left(1 - e^{-p_j t}\right)\right] e^{-(z-1)t} dt, \tag{2.6}$$

or, equivalently, by substituting $x = e^{-t}$ in the integral,

$$G(z) = 1 - (z-1) \int_0^1 \left[1 - \prod_{j=1}^{N} \left(1 - x^{p_j}\right)\right] x^{z-2} dx. \tag{2.7}$$

**Remark 2.1.** *An alternative way to derive (2.6)–(2.7) is by adapting the nice approach of [13], where the main ingredient is an appropriate generating function.*

Observe that,

$$E\left[T_N^{(r)}\right] = (-1)^r \lim_{z \to 1^+} G^{(r)}(z),$$

from which we arrive at the formulas

$$E\left[T_N^{(r)}\right] = r \int_0^\infty \left[1 - \prod_{j=1}^{N} \left(1 - e^{-p_j t}\right)\right] t^{r-1} dt$$

$$= (-1)^{r-1} r \int_0^1 \left[1 - \prod_{j=1}^{N} \left(1 - x^{p_j}\right)\right] \ln(x)^{r-1} \frac{dx}{x}. \tag{2.8}$$

## 2.1 The equally likely case

Naturally, regarding the previous formulas the simplest case occurs when one takes

$$p_1 = \cdots = p_N = \frac{1}{N}. \tag{2.9}$$

Actually, this case apart from its simplicity, has the property that among all sequences, it is the one with the smallest moments of $T_N$. This is a known result (see [5]). For example, (2.7) and (2.8) imply immediately that, for a given $z \geq 1$

$$G(z) = E\left[z^{-T_N}\right],$$

attains its maximum value, while $E\left[T_N^{(r)}\right]$ attain its minimum value, when all $p_j$'s are equal. Under (2.9), one has

$$G(z) = 1 - (z-1) \int_0^1 \left[ 1 - \left(1 - x^{1/N}\right)^N \right] x^{z-2} dx$$

$$= N! \frac{\Gamma\left((z-1)N + 1\right)}{\Gamma(zN+1)},$$

$$E\left[T_N^{(r)}\right] = (-1)^{r-1} r \int_0^1 \left[ 1 - \left(1 - x^{1/N}\right)^N \right] \ln^{r-1}(x) \frac{dx}{x}. \qquad (2.10)$$

Substituting $u = 1 - x^{1/N}$ in the integral of (2.10) one gets

$$E\left[T_N^{(r)}\right] = (-1)^{r-1} r N^r \int_0^1 \frac{1-u^N}{1-u} \ln^{r-1}(1-u) du$$

$$= (-1)^{r-1} r N^r \sum_{m=1}^N \int_0^1 u^{m-1} \ln^{r-1}(1-u) du.$$

Repeated integration by parts in the last integral yields

$$E\left[T_N^{(r)}\right] = r! \, N^r \sum_{m=1}^N \left( \frac{1}{m} \sum_{m_1=1}^m \left( \frac{1}{m_1} \sum_{m_2=1}^{m_1} \frac{1}{m_2} \cdots \right) \sum_{m_{r-1}=1}^{m_{r-2}} \frac{1}{m_{r-1}} \right) = r! N^r \alpha_r(N), \quad (2.11)$$

where the $\alpha_r(N)$'s are defined recursively by

$$\alpha_1(N) = \sum_{m=1}^N \frac{1}{m}, \quad \alpha_r(N) = \sum_{m=1}^N \frac{\alpha_{r-1}(m)}{m}.$$

It seems that formulas for $E\left[T_N^{(r)}\right]$ had been first obtained in [15]. Foata *et al* (see [14]), called the numbers $\alpha_r(N)$ *hyperharmonic* and derived their asymptotics using multivariate generating fuctions. Soon after, Adler *et al* (see [1]), gave explicit expression for the asymptotics of the hyperharmonic numbers using basic probability arguments. In particular, (see [14])

$$\alpha_r(N) \sim \frac{(\ln N)^r}{r!} \quad \text{as} \quad N \to \infty,$$

hence (2.11) yields

$$E\left[T_N^{(r)}\right] \sim N^r (\ln N)^r \quad \text{as} \quad N \to \infty. \qquad (2.12)$$

To give an idea of how higher order asymptotics for $E\left[T_N^{(r)}\right]$ look like, let us mention that, e.g., for $r = 3$ we have either from [14], or by repeated application of Abel partial summation method

$$E\left[T_N^{(3)}\right] = N^3 \left[ \ln^3 N + 3\gamma \ln^2 N + \left(3\gamma^2 + \frac{\pi^2}{2}\right) \ln N \right.$$

$$\left. + \left(2\zeta(3) + \gamma^3 + \frac{\gamma\pi^2}{2}\right) + O\left(\frac{\ln N}{N}\right) \right], \qquad (2.13)$$

where $\gamma$ is the Euler's constant.

## 2.2 Large $N$ asymptotics for general families of coupon probabilities

When $N$ is large it is not obvious at all what information one can obtain for $E\left[T_N^{(r)}\right]$ from formula (2.8). For this reason there is a need to develop efficient ways for deriving asymptotics as $N \to \infty$ (we have already analyzed the very special case of equal probabilities—see formulas (2.12)–(2.13)). Let $\alpha = \{a_j\}_{j=1}^{\infty}$ be a sequence of strictly positive numbers. Then, for each integer $N > 0$, one can create a probability measure $\pi_N = \{p_1, ..., p_N\}$ on the set $\{1, ..., N\}$ by taking

$$p_j = \frac{a_j}{A_N}, \qquad \text{where} \quad A_N = \sum_{j=1}^{N} a_j. \tag{2.14}$$

Notice that $p_j$ depends on $\alpha$ and $N$, thus, given $\alpha$, it makes sense to consider the asymptotic behavior of $E\left[T_N^{(r)}\right]$ as $N \to \infty$. This way of producing sequences of probability measures first appeared in [6].

**Remark 2.2.** *Clearly, for given $N$ the $p_j$'s can be assumed monotone in $j$ without loss of generality. As for the sequence $\{a_j\}_{j=1}^{\infty}$, (i) if $a_j \to \infty$, then for each $k \in \mathbb{N}$ there is a $j = j(k) \geq k$ such that $a_j \geq a_i$, for all $i \leq j$. This tells us that, by rearranging the terms $a_i$, where $j(k) \leq i \leq j(k+1)$, $\{a_j\}_{j=1}^{\infty}$ can be assumed nondecreasing without loss of generality.*
*(ii) Similarly, if $a_j \to 0$, then $\{a_j\}_{j=1}^{\infty}$ can be assumed nonicreasing without loss of generality.*

We set,

$$H_N(\alpha; r) := r \int_0^{\infty} \left[1 - \prod_{j=1}^{N} \left(1 - e^{-a_j t}\right)\right] t^{r-1} \, dt$$

$$= (-1)^{r-1} \, r \int_0^1 \left[1 - \prod_{j=1}^{N} \left(1 - x^{a_j}\right)\right] \ln(x)^{r-1} \frac{dx}{x}. \tag{2.15}$$

If $s\alpha := \{sa_j\}_{j=1}^{\infty}$, by substituting $t = su$ in the first integral of (2.15), we get

$$H_N(s\alpha; r) = s^{-r} H_N(\alpha; r) \tag{2.16}$$

and hence, in view of (2.8) and (2.14),

$$E\left[T_N^{(r)}\right] = A_N^r H_N(\alpha; r). \tag{2.17}$$

As it was noticed in [6] and [8] for $E[T_N]$, the problem of estimating $E\left[T_N^{(r)}\right]$ as $N \to \infty$, can be treated as two separate problems, namely estimating $A_N^r$ and estimating $H_N(\alpha; r)$. Our analysis focuses on estimating $H_N(\alpha; r)$. The estimation of $A_N^r$ will be considered an external matter which can be handled by existing powerful methods, such as the Euler-Maclaurin Summation formula, the Laplace method for sums (see, e.g.,[4]), or even summation by parts.

## 3 Unequal coupon probabilities

### 3.1 The dichotomy

For convenience, we denote

$$f_N^{\alpha}(x) = \prod_{j=1}^{N} (1 - x^{a_j}), \quad 0 \leq x \leq 1.$$

The following properties of the functions $f_N^\alpha$ are immediate:

(i) $f_N^\alpha(0) = 1$ and $f_N^\alpha(1) = 0$, (ii) $f_N^\alpha(x)$ is monotone decreasing in $x$,

(iii) $f_{N+1}^\alpha(x) \le f_N^\alpha(x)$. In particular

$$\lim_N f_N^\alpha(x) = \prod_{j=1}^\infty (1 - x^{a_j}) \quad \text{exists.}$$

Thus, by applying the Monotone Convergence Theorem in (2.15) we get

$$L_r(\alpha) := \lim_N H_N(\alpha; r) = (-1)^{r-1} r \int_0^1 \left[ 1 - \prod_{j=1}^\infty (1 - x^{a_j}) \right] \ln (x)^{r-1} \frac{dx}{x}. \tag{3.1}$$

Notice that $L_r(\alpha) > 0$, for any $\alpha$ (since, for every $x \in (0, 1)$, $f_N^\alpha(x) < 1$ and decreases with $N$). However, we may have $L_r(\alpha) = \infty$. In fact as we will see (in Remark 3.3 below), $L_r(\alpha) = \infty$ if and only if $L_1(\alpha) = \infty$.

**Theorem 3.1.** $L_r(\alpha) < \infty$ if and only if there exist a $\xi \in (0, 1)$ such that

$$\sum_{j=1}^\infty \xi^{a_j} < \infty. \tag{3.2}$$

Before proving the theorem we recall the following lemma (see [6]):

**Lemma 3.2.** Let $\{b_j\}_{j=1}^\infty$ be a sequence of real numbers such that $0 \le b_j \le 1$, for all $j$. If $\sum_{j=1}^\infty b_j < \infty$, then

$$\sum_{j=1}^\infty b_j - \sum_{1 \le l < j} b_l b_j \le 1 - \prod_{j=1}^\infty (1 - b_j) \le \sum_{j=1}^\infty b_j.$$

*Proof of Theorem 3.1.* Assume that there is a $\xi \in (0, 1)$ such that (3.2) is true. Then, by (3.1) and Lemma 3.2 we have that, for all positive integers $r$,

$$L_r(\alpha) \le (-1)^{r-1} r \int_0^\xi \left[ \sum_{j=1}^\infty x^{a_j} \right] \ln (x)^{r-1} \frac{dx}{x}$$

$$+ (-1)^{r-1} r \int_\xi^1 \left[ 1 - \prod_{j=1}^\infty (1 - x^{a_j}) \right] \ln (x)^{r-1} \frac{dx}{x}$$

$$< (-1)^{r-1} r \int_0^\xi \left[ \sum_{j=1}^\infty x^{a_j - 1} \right] \ln (x)^{r-1} dx + (-1)^r r (\ln \xi)^r. \tag{3.3}$$

Now, integration by parts gives

$$I_j(\xi; r) := \int_0^\xi x^{a_j - 1} \ln (x)^{r-1} dx$$

$$= \left[ \frac{x^{a_j} \ln (x)^{r-1}}{a_j} \right]_{x=0}^\xi - (r-1) \int_0^\xi \frac{x^{a_j - 1}}{a_j} \ln (x)^{r-2} dx,$$

hence,

$$I_j(\xi; r) = \frac{1}{a_j} \xi^{a_j} \sum_{k=0}^{r-1} (-1)^k (r-1)_k \frac{1}{a_j^k} \ln (\xi)^{r-1-k}, \tag{3.4}$$

where $(r-1)_k = (r-1)!/(r-1-k)!$ is the falling Pochhammer symbol. Next we apply Fubini-Tonelli's theorem and use (3.4) in (3.3) to get

$$
L_r(\alpha) \leq (-1)^{r-1} r \left( \sum_{j=1}^\infty \frac{1}{a_j} \xi^{a_j} \left( \sum_{k=0}^{r-1} (-1)^k (r-1)_k \frac{1}{a_j^k} \ln^{r-1-k}(\xi) \right) \right)
$$
$$
+ (-1)^r r \ln^r \xi. \tag{3.5}
$$

Now, (3.2) implies that $\xi^{a_j} \to 0$, i.e. $a_j \to \infty$. Therefore, $\min_j \{a_j\} = a_{j_0} > 0$. Thus,

$$
L_r(\alpha) \leq (-1)^{r-1} r \frac{1}{a_{j_0}} \left( \sum_{j=1}^\infty \xi^{a_j} \right) \left( \sum_{k=0}^{r-1} (-1)^k (r-1)_k \frac{1}{a_{j_0}^k} \ln^{r-1-k}(\xi) \right)
$$
$$
+ (-1)^r r \ln^r \xi.
$$

Thus, since $r$ is a positive integer, one obtains $L_r(\alpha) < \infty$ from (3.2) .
Conversely, if $\sum_{j=1}^\infty \xi^{a_j} = \infty$, for all $\xi \in (0,1)$, then, by a well-known property of infinite products (see, e.g. [20])

$$
\prod_{j=1}^\infty (1 - x^{a_j}) = 0, \qquad \text{for all } x \in (0,1)
$$

and hence (3.1) yields $L_r(\alpha) = (-1)^{r-1} r \int_0^1 \left( \ln^{r-1}(x)/x \right) dx = \infty$. $\qquad \square$

**Remark 3.3.** *It has been shown in [6], that $L_1(\alpha) < \infty$, if and only if there exist a $\xi \in (0,1)$ such that $\sum_{j=1}^\infty \xi^{a_j} < \infty$. Thus, $L_r(\alpha) < \infty$ if and only if $L_1(\alpha) < \infty$. To sum up we have the following **dichotomy,** simultaneously for all positive integers $r$:*

$$
\text{(i)} \quad 0 < L_r(\alpha) < \infty \quad \text{or} \quad \text{(ii)} \quad L_r(\alpha) = \infty. \tag{3.6}
$$

**Remark 3.4.** *Consider the error term, defined by*

$$
\Delta_r(N) := L_r(\alpha) - H_N(\alpha; r).
$$

*Then (for all positive integers $r$) by (2.15), (3.1), Fubini-Tonelli's theorem, and repeated integration by parts, we have*

$$
\Delta_r(N) = (-1)^{r-1} r \int_0^1 \prod_{j=1}^N (1 - x^{a_j}) \left[ 1 - \prod_{j=N+1}^\infty (1 - x^{a_j}) \right] \ln(x)^{r-1} \frac{dx}{x}
$$

$$
\leq (-1)^{r-1} r \int_0^1 \left( \sum_{j=N+1}^\infty x^{a_j} \right) \ln(x)^{r-1} \frac{dx}{x} = r! \sum_{j=N+1}^\infty a_j^{-r}. \tag{3.7}
$$

*Thus, if $\sum_{j=1}^\infty a_j^{-r} < \infty$, then (3.7) can serve as an upper bound for the error $\Delta_r(N)$.*

### 3.2 The case $L_r(\alpha)$ **is finite**

Let $A_N$ and $L_r(\alpha)$ be as in (2.14) and (3.1) respectively. We note that, by Theorem 3.1, $L_r(\alpha) < \infty$ implies that $\lim_j a_j = \infty$ (hence $\lim_N A_N = \infty$).

**Theorem 3.5.** *If $L_r(\alpha) < \infty$, then as $N \to \infty$,*

$$
E\left[ T_N^{(r)} \right] = A_N^r L_r(\alpha) \left[ 1 + o(1) \right], \tag{3.8}
$$

*Proof of Theorem 3.5.* Formula (3.8) follows immediately from (2.17) and (3.1). □

Theorem 3.5 states that if $L_r(\alpha) < \infty$, then the asymptotics of $E\left[T_N^{(r)}\right]$ are essentially determined by the asymptotics of $A_N$. As was already mentioned, asymptotic estimates of $A_N$ can be obtained by various known methods. Alternatively, one can resort to specific features of $\alpha$. For instance, if $\alpha$ is of the form

$$a_j = e^{jc_j}, \quad \text{where} \quad c_j \nearrow \infty, \tag{3.9}$$

then it is an easy exercise to see that, as $N \to \infty$,

$$A_N = \sum_{j=1}^{N} a_j \sim a_N. \tag{3.10}$$

To verify (3.10), we use (3.9) and sum a geometric series to get

$$A_N = \sum_{j=1}^{N} e^{jc_j} \leq \sum_{j=1}^{N} e^{jc_N} = \frac{e^{c_N(N+1)} - 1}{e^{c_N} - 1} \leq M e^{c_N(N+1)},$$

where $M = 1/(e^{c_1} - 1)$. Since,

$$\lim_{N \to \infty} \frac{M e^{c_N(N+1)}}{a_{N+1}} = \lim_{N \to \infty} M e^{(c_N - c_{N+1})(N+1)} = 0,$$

the result follows. In words, if a sequence satisfies (3.9), then in the sum of (3.10), the last term dominates all the previous terms. Examples of such sequences are $a_j = e^{j^r}$ with $r > 1$, $a_j = j^j$ and $a_j = j!$ (see Example 4.5).

We now continue with a much more challenging case.

### 3.3 The case $L_r(\alpha)$ is infinite

### 3.3.1 The leading behavior of the rising moments of $T_N$

By Theorem 3.1, $L_r(\alpha) = \infty$ is equivalent to $L_j(\alpha) = \infty$, for all $j = 1, 2, \cdots, r-1$, and also equivalent to $\sum_{j=1}^{\infty} x^{a_j} = \infty$, for all $x \in (0, 1)$. For our further analysis, we follow [6], and write $a_j$ in the form

$$a_j = \frac{1}{f(j)}, \qquad \text{where} \quad f(x) > 0, \tag{3.11}$$

and assume that $f(x)$ possesses two derivatives satisfying the following conditions as $x \to \infty$:

$$\text{(i) } f(x) \nearrow \infty, \quad \text{(ii)} \frac{f'(x)}{f(x)} \searrow 0, \quad \text{and} \quad \text{(iii) } \frac{f''(x)/f'(x)}{[f'(x)/f(x)] \ln [f'(x)/f(x)]} \to 0. \tag{3.12}$$

Conditions (3.12) are satisfied by a variety of commonly used functions. For example,

$$f(x) = x^p (\ln x)^q, \quad p > 0, \ q \in \mathbb{R}, \qquad f(x) = \exp(x^r), \quad 0 < r < 1,$$

as well as various convex combinations of products of such functions.

**Remark 3.6.** *From condition (ii) of (3.12), one has*

$$\lim_{x \to \infty} \frac{f(x+1)}{f(x)} = 1. \tag{3.13}$$

*This can be justified by considering the function $g(x) = \ln(f(x))$ and applying the Mean Value Theorem.*

**Theorem 3.7.** *If* $\alpha = \{1/f(j)\}_{j=1}^{\infty}$, *where* $f$ *satisfies (3.11) and (3.12), then*

$$H_N(\alpha; r) \sim f(N)^r \ln\left(\frac{f(N)}{f'(N)}\right)^r, \qquad N \to \infty. \qquad (3.14)$$

*Proof of Theorem 3.7.* (we adapt the proof of [6] for the leading asymptotics of $H_N(\alpha; 1)$). Set

$$F(x) := -f(x) \ln\left[\frac{f'(x)}{f(x)}\right]. \qquad (3.15)$$

Notice that (3.11) and (ii) of (3.12) imply that $F(x) > 0$, at least for $x$ sufficiently large. Hence, in view of (2.16) one can write (2.15) as:

$$
\begin{aligned}
H_N(\alpha; r) =& F(N)^r H_N\left[\, F(N)\, \alpha;\, r\,\right] \\
=& rF(N)^r \int_0^1 \left[1 - \exp\left(\sum_{j=1}^{N} \ln\left(1 - e^{-\frac{F(N)}{f(j)}s}\right)\right)\right] s^{r-1} ds \\
&+ rF(N)^r \int_1^{\infty} \left[1 - \exp\left(\sum_{j=1}^{N} \ln\left(1 - e^{-\frac{F(N)}{f(j)}s}\right)\right)\right] s^{r-1} ds.
\end{aligned}
\qquad (3.16)
$$

It has been established in [6] that,

$$\lim_N \sum_{j=1}^{N} \ln\left(1 - e^{-\frac{F(N)}{f(j)}s}\right) = \begin{cases} -\infty, & \text{if } s < 1; \\ 0, & \text{if } s \geq 1. \end{cases} \qquad (3.17)$$

and also that

$$\int_1^N e^{-\frac{F(N)}{f(x)}s} dx \sim \frac{1}{s \ln\left[f(N)/f'(N)\right]} \left[\frac{f(N)}{f'(N)}\right]^{1-s}. \qquad (3.18)$$

These two results came out under conditions (3.12). Applying the Bounded Convergence Theorem for the first integral on (3.16) yields (in view of (3.17))

$$
\begin{aligned}
H_N(\alpha; r) =& rF(N)^r \left[\frac{1}{r} + o(1)\right] \\
&+ rF(N)^r \int_1^{\infty} \left[1 - \exp\left(\sum_{j=1}^{N} \ln\left(1 - e^{-\frac{F(N)}{f(j)}s}\right)\right)\right] s^{r-1} ds.
\end{aligned}
\qquad (3.19)
$$

Next, we want to estimate the integral which appears in (3.19). We begin by noticing that by the Dominated Convergence Theorem (since $f(N)/f'(N) \to \infty$)

$$\lim_N \int_1^{\infty} \left[1 - \exp\left(-\frac{(f(N)/f'(N))^{1-s}}{s \ln(f(N)/f'(N))}\right)\right] s^{r-1} ds = 0.$$

In view of (3.18) the above formula implies that

$$\lim_N \int_1^{\infty} \left[1 - \exp\left(-\int_1^N e^{-\frac{F(N)}{f(x)}s} dx\right)\right] s^{r-1} ds = 0. \qquad (3.20)$$

Since $f$ is increasing, we have

$$
\begin{aligned}
\int_1^N e^{-\frac{F(N)}{f(x)}s} dx &\leq \sum_{j=1}^{N} e^{-\frac{F(N)}{f(j)}s} \\
&\leq \int_1^{N+1} e^{-\frac{F(N)}{f(x)}s} dx \\
&\leq \int_1^N e^{-\frac{F(N)}{f(x)}s} dx + e^{-\frac{F(N)}{f(N+1)}s}.
\end{aligned}
\qquad (3.21)
$$

From the above inequalities it follows

$$1 - \exp\left(-\int_1^N e^{-\frac{F(N)}{f(x)}s}dx\right) \leq 1 - \exp\left(-\sum_{j=1}^N e^{-\frac{F(N)}{f(j)}s}\right)$$

$$\leq 1 - \exp\left(-\int_1^N e^{-\frac{F(N)}{f(x)}s}dx + e^{-\frac{F(N)}{f(N+1)}s}\right). \qquad (3.22)$$

However, by (3.18)

$$\lim_N \int_1^N e^{-\frac{F(N)}{f(x)}s}dx = \begin{cases} \infty, & \text{if } s < 1; \\ 0, & \text{if } s \geq 1. \end{cases} \qquad (3.23)$$

Hence, by taking limits in (3.22) and using (3.20) and (3.13), we get

$$\lim_N \int_1^\infty \left[1 - \exp\left(\sum_{j=1}^N \ln\left(1 - e^{-\frac{F(N)}{f(j)}s}\right)\right)\right] s^{r-1}ds = 0. \qquad (3.24)$$

Finally, by the definition of $F(\cdot)$ and the Taylor expansion for the logarithm, namely $\ln(1-x) \sim -x$ as $x \to 0$, (3.19) yields

$$H_N(\alpha;r) \sim F(N)^r = f(N)^r \ln\left(\frac{f(N)}{f'(N)}\right)^r, \qquad N \to \infty \qquad (3.25)$$

and the proof is completed. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\square$

**Remark 3.8.** *Using Theorem 3.7 in (2.17) we get, as $N \to \infty$,*

$$E\left[T_N^{(r)}\right] \sim A_N^r f(N)^r \ln\left(\frac{f(N)}{f'(N)}\right)^r = \frac{1}{\min_{1 \leq j \leq N}\{p_j\}^r} \ln\left(\frac{f(N)}{f'(N)}\right)^r, \qquad (3.26)$$

*where the last equality follows from (2.14).*

### 3.3.2 Asymptotic estimates for the rising moments of $T_N$ by comparison with known sequences

In this subsubsection we will present a theorem that helps us obtain asymptotic estimates by comparison with sequences $\alpha$ for which the asymptotic estimates of $H_N(\alpha;r)$ are known (for instance, via Theorem 3.7). A similar theorem concerning the special case of $r = 1$, can be found in [6]. First, we recall the following notation. Suppose that $\{s_j\}_{j=1}^\infty$ and $\{t_j\}_{j=1}^\infty$ are two sequences of nonnegative terms. The symbol $s_j \asymp t_j$ means that there are two constants $C_1 > C_2 > 0$ and an integer $j_0 > 0$ such that

$$C_2 t_j \leq s_j \leq C_1 t_j, \qquad \text{for all } j \geq j_0, \qquad (3.27)$$

i.e. $s_j = O(t_j)$ and $t_j = O(s_j)$.

**Theorem 3.9.** *Let $\alpha = \{a_j\}_{j=1}^\infty$ and $\beta = \{b_j\}_{j=1}^\infty$ be sequences of strictly positive terms such that $\lim_N H_N(\alpha;r) = \lim_N H_N(\beta;r) = \infty$.*
*(i) If there exists an $j_0$ such that $a_j = b_j$, for all $j \geq j_0$, then $H_N(\beta;r) - H_N(\alpha;r)$ is bounded,*
*(ii) if $a_j = O(b_j)$, then $H_N(\beta;r) = O(H_N(\alpha;r))$ as $N \to \infty$,*
*(iii) if $a_j = o(b_j)$, then $H_N(\beta,r) = o(H_N(\alpha;r))$ as $N \to \infty$,*
*(iv) if $a_j \asymp b_j$, then $H_N(\beta;r) \asymp H_N(\alpha;r)$ as $N \to \infty$,*
*(v) if $a_j \sim b_j$, then $H_N(\beta;r) \sim H_N(\alpha;r)$ as $N \to \infty$.*

*Proof of Theorem 3.9.* We will prove (i) and (v). The proofs of (ii)–(iv) are similar. Case (i) follows easily from (2.15):

$$|H_N(\beta;r) - H_N(\alpha;r)|$$

$$= r\left|\int_0^\infty \prod_{j=j_0}^N \left(1 - e^{-a_j t}\right)\left[\prod_{j=1}^{j_0-1}\left(1 - e^{-a_j t}\right) - \prod_{j=1}^{j_0-1}\left(1 - e^{-b_j t}\right)\right] t^{r-1}\, dt\right|$$

$$\leq \int_0^\infty \left|\left[\prod_{j=1}^{j_0-1}\left(1 - e^{-a_j t}\right) - \prod_{j=1}^{j_0-1}\left(1 - e^{-b_j t}\right)\right]\right| t^{r-1}\, dt$$

$$= \int_0^\infty \left|\sum_{J\subset\{1,...,j_0-1\}} (-1)^{|J|}\left\{\exp\left(-t\sum_{j\in J} a_j\right) - \exp\left(-t\sum_{j\in J} b_j\right)\right\} t^{r-1}\right| dt\ <\infty,$$

where we have used (2.4). The sum extends over all $2^{j_0-1}$ subsets $J$ of $\{1,...,j_0-1\}$, while $|J|$ denotes the cardinality of $J$.

To prove (v) we first fix an $\epsilon > 0$. Then $(1-\epsilon)b_j \leq a_j \leq (1+\epsilon)b_j$, for all $j \geq j_0(\epsilon)$. Thus, by case (i) there is an $M = M(\epsilon)$ such that

$$\left(\frac{1}{1+\epsilon}\right)^r H_N(\beta;r) - M \leq H_N(\alpha;r) \leq \left(\frac{1}{1-\epsilon}\right)^r H_N(\beta;r) + M,$$

for all $N \geq N_0(\epsilon)$. If we divide by $H_N(\beta;r)$ and then let $N \to \infty$, we obtain (v) since $\epsilon$ is arbitrary and $\lim_N H_N(\beta;r) = \infty$. □

## 4 Examples

**Example 4.1.** *The case $a_j = 1$, for all $j$, has been already discussed in detail in Section 2. This case can also provide us with an application of Theorem 3.9: If $\beta = \{b_j\}_{j=1}^\infty$ is a sequence such that $0 < \underline{\lim} b_j \leq \overline{\lim} b_j < \infty$ then, there are two constants $C_1 > C_2 > 0$ and an integer $j_0 > 0$ such that*

$$C_2 b_j \leq 1 \leq C_1 b_j, \qquad \text{for all } j \geq j_0, \quad i.e. \ \ 1 \asymp b_j.$$

*Hence, by part (iv) of Theorem 3.9, $H_N(\beta;r) \asymp \ln^r N$. If, in addition, $\lim b_j = b$ exists, then $ba_j \sim b_j$. Hence, by part (v) of Theorem 3.9, $H_N(\beta;r) \sim H_N(b\alpha;r)$. Using (2.16) we get*

$$H_N(\beta;r) \sim b^{-r}\ln^r N.$$

**Example 4.2.** *$a_j = j^p$, where $p > 0$. In this case*

$$L_{r,p} := L_r(\alpha) = (-1)^{r-1}\, r\int_0^1\left[1 - \prod_{j=1}^\infty (1 - x^{j^p})\right]\ln^{r-1}(x)\frac{dx}{x},$$

*(notice that $L_{r,p}$ decrease with $p$). By Theorem 3.1 and for all positive integers $r$ we have: $L_{r,p} < \infty$. Now, in accordance with (3.8) we also need to estimate $A_N$. From the Euler-Maclaurin summation formula we get the full asymptotic expansion of $A_N = \sum_{n=1}^N n^p$ (in fact, if $p$ is a positive integer, $A_N$ is a polynomial in $N$ of degree $p+1$). In particular,*

$$A_N = \sum_{n=1}^N n^p = \frac{N^{p+1}}{p+1}\left[1 + O\left(\frac{1}{N}\right)\right].$$

*Therefore, by (2.17)*

$$E\left[T_N^{(r)}\right] = \frac{N^{r(p+1)}}{(p+1)^r} L_{r,p}\left[1 + o(1)\right].$$

*The case $p = 1$ is known as the linear case, and it is of particular interest. From Euler's pentagonal-number formula (a combinatorial proof by F. Franklin can be found, e.g., in [2])*

$$\prod_{j=1}^{\infty}\left(1 - x^j\right) = 1 + \sum_{k=1}^{\infty}(-1)^k\left[x^{\omega(k)} + x^{\omega(-k)}\right],$$

$$\omega(k) = (3k^2 - k)/2, \quad k = 0, \pm 1, \pm 2, \dots$$

*In that case $L_r$ becomes*

$$L_r = (-1)^r\, r \sum_{k=1}^{\infty}(-1)^k\left[\int_0^1 x^{\omega(k)-1}\ln(x)^{r-1}dx + \int_0^1 x^{\omega(-k)-1}\ln(x)^{r-1}dx\right].$$

*Repeated integration by parts yields,*

$$L_r = r!\sum_{k=1}^{\infty}(-1)^{k+1}\left\{\frac{1}{\omega(k)^r} + \frac{1}{\omega(-k)^r}\right\}$$

$$= 2^r\, r!\sum_{k=1}^{\infty}(-1)^{k+1}\left[\frac{1}{(3k^2 - k)^r} + \frac{1}{(3k^2 + k)^r}\right].$$

*For example, (see [6], [8])*

$$L_1 = \frac{4\pi\sqrt{3}}{3} - 6 \cong 1.2552, \quad L_2 = 4(54 - 8\pi\sqrt{3} - \pi^2) \cong 2.39684.$$

*As for $L_3$, a numerical computation gives*

$$L_3 \cong 6.68903.$$

**Example 4.3.** $b_j = e^{pj}$, $a_j = e^{-pj}$, $p > 0$. *For the sequence $\beta = \{b_j\}_{j=0}^{\infty}$ we have, $L_r(\beta) < \infty$, $r = 1, 2, \cdots$. Furthermore,*

$$\Delta_r(N) = L_r(\beta) - H_N(\beta; r) \leq r!\sum_{j=N+1}^{\infty} e^{-rpj} = \frac{r!\, e^{-rp(N+1)}}{1 - e^{-rp}}, \qquad (4.1)$$

$$B_N := \sum_{j=0}^{N} b_j = \frac{e^{p(N+1)} - 1}{e^p - 1}, \quad \text{and} \quad E\left[T_N^{(r)}\right] \sim \left(\frac{e^{p(N+1)}}{e^p - 1}\right)^r L_r(\beta).$$

*In the special case of $b_j = 2^j$ (i.e. $p = \ln 2$), we have*

$$\phi(x) := \prod_{j=0}^{\infty}(1 - x^{2^j}) = \sum_{k=0}^{\infty}(-1)^{\delta(k)}x^k, \qquad (4.2)$$

*where $\delta(k)$ is the number of ones in the binary expansion of $k$. Therefore, by (3.1)*

$$L_r(\beta) = (-1)^r\, r\sum_{k=1}^{\infty}(-1)^{\delta(k)}\left[\int_0^1 x^{k-1}\ln(x)^{r-1}dx\right] = r!\sum_{k=1}^{\infty}\frac{(-1)^{\delta(k)-1}}{k^r}. \qquad (4.3)$$

*Now, for the sequence $\alpha = \{a_j\}_{j=0}^{\infty}$ we have $L_r(\alpha) = \infty$. Furthermore $f(x) = e^{px}$ does not satisfy condition (ii) and of (3.12), thus Theorem 3.7 cannot be applied. However,*

*if we let $c_N = e^{pN}$, then $\{b_j : 0 \le j \le N\} = \{c_N a_j : 0 \le j \le N\}$, for each $N$, i.e. the elements of the two truncated sequences are proportional to each other. Hence, the sequences $\beta$ and $\alpha$ produce the same coupon probabilities. In this way we get cheaply a counterexample for Theorem 3.7, in case where $f(\cdot)$ does not satisfy all conditions of (3.12).*

**Example 4.4.** *$a_j = 1/j^p$, $p > 0$. This is the so-called generalized Zipf law. In this case Theorem 3.1 implies $L_r(\alpha) = \infty$. If $f(x) = x^p$, then $f$ satisfies (3.12) and hence Theorem 3.7 apply. It is now straightforward (say, form the Euler-Maclaurin Summation formula—see, e.g., [2]) to estimate $A_N^r$ and get*

$$
\begin{aligned}
&A_N^r \sim \left(\frac{1}{1-p}\right)^r N^{r-rp}, && \text{if} \quad 0 < p < 1, \\
&A_N^r = H_N^r \sim \ln^r N, && \text{if} \quad p > 1, \\
&A_N^r \sim \zeta(p)^r, && \text{if} \quad p > 1,
\end{aligned}
$$

*where $\zeta(\cdot)$ denotes the Riemann zeta function. Hence Theorem 3.7 gives*

$$
\begin{aligned}
&E\left[T_N^{(r)}\right] \sim \frac{N^r \ln^r(N)}{(1-p)^r} && \text{for} \quad 0 < p < 1, \\
&E\left[T_N^{(r)}\right] \sim N^r \ln^{2r}(N) && \text{for} \quad p = 1, \\
&E\left[T_N^{(r)}\right] \sim \zeta(p)^r N^{rp} \ln^r(N) && \text{for} \quad p > 1.
\end{aligned}
$$

**Example 4.5.** *$a_j = j!$. We have $L_r(\alpha) < \infty$. Here, the Euler-Maclaurin summation formula is not effective for the estimation of $A_N$. However, Stirling's formula and (3.9)–(3.10) imply easily that*

$$A_N \sim N!.$$

*Hence, by Theorem 3.5 we get*

$$E\left[T_N^{(r)}\right] \sim L_r(\alpha)(N!)^r \quad \text{as} \quad N \to \infty.$$

## 5 Concluding remarks

The main topic of this paper was the asymptotics of $E[T_N^{(r)}]$, namely the $r$-th rising moment of $T_N$, as $N \to \infty$. We have already mentioned the work of H.J. Godwin [15], in the case of uniform coupon probabilities. We are not aware of any previous work on asymptotics of higher rising moments ($r \ge 3$) in the case of unequal coupon probabilities. Of course, in the existing literature there are many works on the asymptotics of $E[T_N]$ and, also, few works regarding $E[T_N^2]$ and $V[T_N]$, namely the variance of the random variable $T_N$.

Let us discuss briefly few representative works. The first and the second moment of $T_N$ were studied in [7]. In this article R.K. Brayton (Ph.D. thesis under N. Levinson) derived an asymptotic formula for $V[T_N]$ under very restrictive assumptions on $\alpha$. In particular, the probabilities $p_j$ considered in [7] must satisfy:

$$
\lambda(N) := \frac{\max_{1 \le j \le N}\{p_j\}}{\min_{1 \le j \le N}\{p_j\}} \le M < \infty, \qquad \text{independently of } N. \tag{5.1}
$$

General asymptotic estimates, for the case $r = 1$ were found in [6], for the families of coupon probabilities which we study in the present paper. Our results here are in accordance with [6].

The case $r = 1, 2$ was considered in [8]. The authors adopted the dichotomy of [6] and obtained the leading behavior of the variance $V[T_N]$. Moreover, for a large class of families of coupon probabilities they obtained detailed asymptotics of $E[T_N]$ and $E[T_N(T_N + 1)]$ (up to the fifth and sixth term respectively). Notice that their results complement the results of [7], since they concern quite general sequences for which the ratio $\lambda(N)$ of (5.1) is not bounded (e.g. linear and Zipf).

Recently, J. Du Boisberranger, D. Gardy, and Y. Ponty, [9] considered the word collector problem, i.e. the expected number of calls to a random weighted generator before all the words of a given length in a language are generated. The main ingredient of this instance of the non-uniform coupon collector lies in the, potentially large, multiplicity of the words (coupons) of a given probability (composition). They obtained a general theorem that gives an asymptotic equivalent for the expected waiting time of a general version of the Coupon Collector (case $r = 1$). This theorem is especially well-suited for classes of coupons featuring high multiplicities. Their results and [6] are complementary.

Finally, let us mention that it was pointed out to us that an important case in the applications is when there is a subcollection of coupons that continues to grow (as $N$ grows) and all of the coupons in the subcollection have the same probability; it could be called "uniform subcollection". This can be modeled by a sequence $\{a_j\}_{j=1}^{\infty}$ which is the "union" of two subsequences one of which is constant (this corresponds to the uniform subcollection of coupons), while the other is of the form discussed in this subsection. In this case, we conjecture that Theorem 3.7 is still valid (under an appropriate renaming of the index) provided the "density" of the constant subsequence is sufficiently small. In other words, the uniform subcollection does not affect the asymptotic distribution.

Furthermore, if $\{a_j\}_{j=1}^{\infty}$ is the "union" of two vanishing subsequences one of which decays much faster than the other, then we conjecture that the faster one prevails in the asymptotics, provided its density is not very small.

# References

[1] Adler, I., Oren, S., and Ross, S.: The coupon collector's problem revisited, *J. Appl. Prob.* **40** (2003), no. 2, 513–518. MR-1978107

[2] Apostol, T. M.: *Introduction to Analytic Number Theory*, Springer-Verlag, New York-Heidelberg, 1976. xii+338 pp. MR-0434929

[3] Baum, L. E. and Billingsley, P.: Asymptotic distributions for the coupon collector's problem, *Ann. Math. Statist.* **36** (1965) 1835–1839. MR-0182039

[4] Bender, C. M. and Orszag, S. A.: *Advanced Mathematical Methods for Scientists and Engineers I: Asymptotic Methods and Perturbation Theory*, Springer-Verlag, New York, 1999. xiv+593 pp. MR-1721985

[5] Boneh, A. and Hofri, M.: The coupon collector problem revisited–A survey of engineering problems and computational methods, *Comm. Statist. Stochastic Models* **13** (1997), no. 1, 39–66. MR-1430927

[6] Boneh, S. and Papanicolaou, V. G.: General asymptotic estimates for the coupon collector problem, *J. Comput. Appl. Math.* **67** (1996), no. 2, 277–289. MR-1390185

[7] Brayton, R. K.: On the asymptotic behavior of the number of trials necessary to complete a set with random selection, *J. Math. Anal. Appl.* **7**(1963) 31–61. MR-0158427

[8] Doumas, A. V. and Papanicolaou, V. G.: The coupon collector's problem revisited: asymptotics of the variance, *Adv. in Appl. Probab.* **44** (2012), no. 1, 166–195. MR-2951551

[9] Du Boisberranger, J., Gardy, D. and Ponty, Y.: The weighted words collector, *23rd Intern. Meeting on Probabilistic, Combinatorial, and Asymptotic Methods for the Analysis of Algorithms (AofA'12)*, 243–264, Discrete Math. Theor. Comput. Sci. Proc., AQ, Assoc. Discrete Math. Theor. Comput. Sci., Nancy, 2012. MR-2957335

[10] Durrett, R.: *Probability: theory and examples*, Second edition. Duxbury Press, Belmont, CA, 1996. xiii+503 pp. MR-1609153

[11] Erdős, P. and Rényi, A.: On a classical problem of probability theory, *Magyar. Tud. Akad. Mat. Kutató Int. Kőzl.*, **6** (1961), 215–220. MR-0150807

[12] Feller, W.: *An Introduction to Probability Theory and Its Applications, Vol. I*, Third Edition, John Wiley & Sons, Inc., New York-London-Sydney 1968 xviii+509 pp. MR-0228020

[13] Flajolet, P., Gardy, D. and Thimonier, L.: Birthday paradox, coupon collectors, caching algorithms and self-organizing search, *Discrete Appl. Math.* **39** (1992), no. 3, 207–229. MR-1189469

[14] Foata, D., Guo-Niu, H. and Lass, B.: Les nombres hyperharmonique et la fratrie du collectionneur de vignettes.(French) [Hyperharmonic numbers and the coupon collector's brotherhood] *Se'm. Lothar. Combin.* **47** (2001/02), Article B47a, 20 pp. MR-1894021

[15] Godwin, H. J.: On cartophily and motor cars, *Math. Gazette* **33** (1949) 169–171.

[16] Holst, L.: On birthday, collectors', occupancy and other classical urn problems, *Internat. Statist. Rev.***54** (1986), no. 1, 15–27. MR-0959649

[17] Janson, S.: Limit theorems for some sequential occupancy problems, *J. Appl. Probab.***20** (1983), no. 3, 545–553. MR-0713504

[18] Neal, P.: The generalized coupon collector problem, *J. Appl. Prob.* **45** (2008), no. 3, 621–629. MR-2455173

[19] Newman, D. J. and Shepp, L.: The double Dixie cup problem, *Amer. Math. Monthly* **67** (1960) 58–61. MR-0120672

[20] Rudin, W.: *Real and Complex Analysis*, Third edition. McGraw-Hill Book Co., New York, 1987. xiv+416 pp. MR-0924157