



93.16 the Coupon Collector's Problem with Unequal Probabilities

Author(s): Stella Dudzic

Source: *The Mathematical Gazette*, Vol. 93, No. 526 (Mar., 2009), pp. 126-130

Published by: The Mathematical Association

Stable URL: <http://www.jstor.org/stable/40378689>

Accessed: 27-07-2016 07:29 UTC

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at
<http://about.jstor.org/terms>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



The Mathematical Association is collaborating with JSTOR to digitize, preserve and extend access to *The Mathematical Gazette*

93.16 The coupon collector's problem with unequal probabilities

The classic coupon collector's problem, where k different items are to be collected (say by opening packets with one coupon in each) with each item having probability $\frac{1}{k}$ of occurring, is well known (see [1]). The expected number of packets which need to be opened to give a full set is

$$k \sum_{i=1}^k \frac{1}{i}. \quad (1)$$

The unequal probabilities case recently came up in a *Times Educational Supplement* online forum [2]. The problem was solved in 1954 by Herman Von Schelling [3]; however, the following solution uses only techniques available to A Level students and results in a different form for the formula for the expected number of packets. Comparison with the equal probability case leads to a result concerning the harmonic series.

Two kinds of coupon

Consider two types of coupon, A and B , with probabilities p and q , respectively, where $p + q = 1$.

Consider buying n packets and opening them all, one at a time. Let $P_{n,2}$ be the probability that there are at least one of each of type A and type B in n packets, that is the probability of completing a full set in n packets, or fewer.

The terms of $(A + B)^n$ give all the different combinations of A and B in the n packets and the corresponding terms of $(p + q)^n$ give the probabilities of these combinations.

$$P_{n,2} = (p + q)^n - p^n - q^n. \quad (2)$$

(This simplifies to $P_{n,2} = 1 - p^n - q^n$, which is the same as saying that getting at least one of each is the complement of them all being the same, but the unsimplified formula holds even when $p + q < 1$.)

Three kinds of coupon

Consider three types of coupon, A , B and C with probabilities p , q and r , respectively, where $p + q + r = 1$.

Consider buying n packets and opening them all, one at a time. Let $P_{n,3}$ be the probability that there are at least one of each of type in n packets, that is the probability of completing a full set in n packets, or fewer.

The terms of $(A + B + C)^n$ give all the different combinations of A , B and C in the n packets and the corresponding terms of $(p + q + r)^n$ give the probabilities of these combinations.

For at least one of each type of coupon, we need to subtract the probabilities of getting exactly two kinds and of getting all n the same. Thus

$$\begin{aligned} P_{n,3} &= (p+q+r)^n - \{(p+q)^n - p^n - q^n\} - \{(q+r)^n - q^n - r^n\} \\ &\quad - \{(r+p)^n - r^n - p^n\} - p^n - q^n - r^n \\ &= (p+q+r)^n - \{(p+q)^n + (q+r)^n + (r+p)^n\} + \{p^n + q^n + r^n\}. \quad (3) \end{aligned}$$

This formula still holds if $p+q+r < 1$ since it consists of those terms of $(p+q+r)^n$ that include all three of p, q, r .

Working inductively

Similarly, for four kinds of coupon with probabilities p, q, r and s

$$P_{n,4} = (p+q+r+s)^n - \sum_4 (p+q+r)^n + \sum_4 (p+q)^n - \sum_4 p^n \quad (4)$$

where \sum_4 denotes summation of all such terms which are combinations from the four probabilities p, q, r and s . This formula can be derived by writing out all the terms and collecting like terms, but the amount of writing this requires becomes increasingly daunting for larger numbers of coupons. Proceeding as shown below forms the outline of a method that can be used to prove the general case, by induction.

$$P_{n,4} = (p+q+r+s)^n - \{\sum_4 P_{n,3} + \sum_4 P_{n,2} + \sum_4 p^n\}$$

$$\sum_4 P_{n,3} = \sum_4 (p+q+r)^n - \sum_4 \{\sum_3 (p+q)^n\} + \sum_4 \{\sum_3 p^n\}.$$

Each $\sum_3 (p+q)^n$ contains 3C_2 terms like $(p+q)^n$ hence $\sum_4 \{\sum_3 (p+q)^n\}$ contains ${}^4C_3 \times {}^3C_2$ such terms. Moreover, $\sum_4 \{\sum_3 (p+q)^n\}$ is a symmetric function and so it is a multiple of $\sum_4 (p+q)^n$, which has 4C_2 terms.

Therefore $\sum_4 \{\sum_3 (p+q)^n\} = \frac{{}^4C_3 \times {}^3C_2}{{}^4C_2} \sum_4 (p+q)^n$.

The terms to be subtracted from $(p+q+r+s)^n$ are:

$$\sum_4 P_{n,3} = \sum_4 (p+q+r)^n - \frac{{}^4C_3 \times {}^3C_2}{{}^4C_2} \sum_4 (p+q)^n + \frac{{}^4C_3 \times {}^3C_1}{{}^4C_1} \sum_4 p^n$$

$$\sum_4 P_{n,2} = \sum_4 (p+q)^n - \frac{{}^4C_2 \times {}^2C_1}{{}^4C_1} \sum_4 p^n$$

$$\sum_4 p^n = \sum_4 p^n$$

and the total of these terms is $\sum_4 (p+q+r)^n - \sum_4 (p+q)^n + \sum_4 p^n$, which gives formula (4).

Binomial coefficients

Expressions of the form $\frac{{}^aC_b \times {}^bC_d}{{}^aC_d}$ simplify.

$$\frac{{}^aC_b \times {}^bC_d}{{}^aC_d} = \frac{a!}{b!(a-b)!} \frac{b!}{d!(b-d)!} \frac{d!(a-d)!}{a!} = \frac{(a-d)!}{(a-b)!(b-d)!} = {}^{a-d}C_{a-b}.$$

Five kinds of coupon

For five kinds of coupon with probabilities p, q, r, s and t

$$P_{n,5} = (p + q + r + s + t)^n - \sum_5 (p + q + r + s)^n + \sum_5 (p + q + r)^n - \sum_5 (p + q)^n + \sum_5 p^n. \quad (5)$$

The terms to be subtracted from $(p + q + r + s + t)^n$ are:

$$\begin{aligned} \sum_5 P_{n,4} &= \sum_5 (p + q + r + s)^n - \frac{{}^5C_4 \times {}^4C_3}{{}^5C_3} \sum_5 (p + q + r)^n + \frac{{}^5C_4 \times {}^4C_2}{{}^5C_2} \sum_5 (p + q)^n - \frac{{}^5C_4 \times {}^4C_1}{{}^5C_1} \sum_5 p^n \\ \sum_5 P_{n,3} &= \sum_5 (p + q + r)^n - \frac{{}^5C_3 \times {}^3C_2}{{}^5C_2} \sum_5 (p + q)^n + \frac{{}^5C_3 \times {}^3C_1}{{}^5C_1} \sum_5 p^n \\ \sum_5 P_{n,2} &= \sum_5 (p + q)^n - \frac{{}^5C_2 \times {}^2C_1}{{}^5C_1} \sum_5 p^n \\ \sum_5 p^n &= \sum_5 p^n \end{aligned}$$

Adding these four and simplifying binomial expressions gives

$$\begin{aligned} &\sum_5 (p + q + r + s)^n + (-{}^2C_1 + 1) \sum_5 (p + q + r)^n \\ &+ ({}^3C_1 - {}^3C_2 + 1) \sum_5 (p + q)^n + (-{}^4C_1 + {}^4C_2 - {}^4C_3 + 1) \sum_5 p^n. \end{aligned}$$

This simplifies to

$$\begin{aligned} &\sum_5 (p + q + r + s)^n + \{(1 - 1)^2 - 1\} \sum_5 (p + q + r)^n \\ &+ \{(-1 + 1)^3 + 1\} \sum_5 (p + q)^n + \{(1 - 1)^4 - 1\} \sum_5 p^n, \end{aligned}$$

that is $\sum_5 (p + q + r + s)^n - \sum_5 (p + q + r)^n + \sum_5 (p + q)^n - \sum_5 p^n$ and subtracting from $(p + q + r + s + t)^n$ gives (5).

The probability of completing the set on opening the n th packet

When collecting k coupons, the probability of completing the full set on opening the n th packet is

$$\begin{aligned} P_{n,k} - P_{n-1,k} &\text{ for } n \geq k; \\ 0 &\text{ otherwise.} \end{aligned} \quad (6)$$

Expected number of packets needed to complete a full set of coupons

Let X be the number of packets needed to complete a full set of coupons.

For $k = 5$ using P_r to represent $P_{r,5}$, we have $E(X) = \lim_{n \rightarrow \infty} E_n(X)$ where

$$\begin{aligned}
 E_n(X) = & \quad nP_n \quad -nP_{n-1} \\
 & + (n-1)P_{n-1} \quad - (n-1)P_{n-2} \\
 & + (n-2)P_{n-2} \quad - (n-2)P_{n-3} \\
 & \quad \quad \quad \vdots \quad \quad \quad \vdots \\
 & \quad \quad \quad + 7P_7 \quad \quad \quad - 7P_6 \\
 & \quad \quad \quad + 6P_6 \quad \quad \quad - 6P_5 \\
 & \quad \quad \quad + 5P_5
 \end{aligned}$$

and therefore

$$\begin{aligned}
 E_n(X) &= nP_n - \sum_{i=5}^{n-1} P_i \\
 &= nP_n - \sum_{i=5}^{n-1} \{ (p+q+r+s+t)^i - \sum_5 (p+q+r+s)^i + \sum_5 (p+q+r)^i \\
 &\quad - \sum_5 (p+q)^i + \sum_5 p^i \}.
 \end{aligned}$$

$$\text{Now } \sum_{i=5}^{n-1} (p+q+r+s+t)^i = \sum_{i=5}^{n-1} 1 = n-5 \text{ since } p+q+r+s+t=1$$

and

$$\begin{aligned}
 \sum_{i=5}^{n-1} (p+q+r)^i &= \frac{(p+q+r)^5 \{1 - (p+q+r)^{n-5}\}}{\{1 - (p+q+r)\}} \\
 &= \frac{(p+q+r)^5 \{1 - (p+q+r)^{n-5}\}}{(s+t)}.
 \end{aligned}$$

Similar results hold for all the other geometric series in $E_n(X)$.

Thus $E_n(X) = nP_n - (n-5) +$ sums of geometric series which have finite sums to infinity. As $n \rightarrow \infty$, $P_n \rightarrow 1$, hence

$$\begin{aligned}
 E(X) &= 5 + \sum_5 \frac{(p+q+r+s)^5}{t} - \sum_5 \frac{(p+q+r)^5}{(s+t)} \\
 &\quad + \sum_5 \frac{(p+q)^5}{(r+s+t)} - \sum_5 \frac{p^5}{(q+r+s+t)}. \tag{7}
 \end{aligned}$$

With similar results for values of k other than 5.

The equal probability case

If $p = q = r = s = t = \frac{1}{5}$, then (7) becomes

$$E(X) = 5 + {}^5C_4 \left(\frac{4}{5}\right)^5 \left(\frac{5}{1}\right) - {}^5C_3 \left(\frac{3}{5}\right)^5 \left(\frac{5}{2}\right) + {}^5C_2 \left(\frac{2}{5}\right)^5 \left(\frac{5}{3}\right) - {}^5C_1 \left(\frac{1}{5}\right)^5 \left(\frac{5}{4}\right).$$

This has to be give the same as the answer given by (1), hence

$$5 \sum_{i=1}^5 \frac{1}{i} = 5 + {}^5C_4 \left(\frac{4}{5}\right)^5 \left(\frac{5}{1}\right) - {}^5C_3 \left(\frac{3}{5}\right)^5 \left(\frac{5}{2}\right) + {}^5C_2 \left(\frac{2}{5}\right)^5 \left(\frac{5}{3}\right) - {}^5C_1 \left(\frac{1}{5}\right)^5 \left(\frac{5}{4}\right).$$

More generally, for $n \geq 2$,

$$n \sum_{i=1}^n \frac{1}{i} = n + \sum_{i=1}^{n-1} {}^nC_i (-1)^{i+1} \left(\frac{n-i}{n}\right)^n \left(\frac{n}{i}\right).$$

References

1. F. Mosteller, *Fifty challenging problems in probability with solutions*, Dover (1965).
2. <http://community.tes.co.uk/forums/t/110234.aspx?PageIndex=1>
3. Herman Von Schelling, Coupon collecting for unequal probabilities, *The Amer. Math. Monthly*, Vol. 61, No. 5. (May, 1954), pp. 306-311.

STELLA DUDZIC

MEI, Monckton House, Epsom Centre, White Horse Business Park,
Trowbridge BA14 0XG

93.17 Head-tail imbalance

When studying probability, one of the first situations students study is that of repeatedly throwing an unbiased coin and counting the number of heads that appear. It is well known that there are approximately the same number of heads and tails, and the result of such an experiment shows that we are unlikely to get exactly the same number of heads as tails. We are going to consider that imbalance in this article, and we shall consider the following three questions in particular.

1. For a particular number of throws, what is the distribution of the numerical difference between the numbers of heads and of tails?
2. If we throw until we have a given number of either heads or tails, how many of the other outcome will we have?
3. How long will it take for the difference in frequencies to exceed a given value?

Throughout, we shall let N stand for the total number of throws and p, q the probabilities of a head and a tail respectively. For much of the work $p = q = \frac{1}{2}$.

Question 1

As is well known, when a coin is thrown N times the number of heads X is given by the binomial distribution, so that:

$$P(X = r) = {}^NC_r p^r q^{N-r}.$$