# Counting sub-multisets of fixed cardinality

Sebastiano Ferraris[*1], Alex Mendelson[†1], Gerardo Ballesio[‡2] and
Tom Vercauteren[§1]

[1]Translational Imaging Group, Centre for Medical Image Computing,
Wolfson House, 4 Stephenson Way, University College London.
[2]Independent researcher.

November 20, 2015

## Abstract

This report presents an expression for the number of a multiset's sub-multisets of a given cardinality as a function of the multiplicity of its elements. This is also the number of distinct samples of a given size that may be produced by sampling without replacement from a finite population partitioned into subsets, in the case where items belonging to the same subset are considered indistinguishable. Despite the generality of this problem, we have been unable to find this result published elsewhere.

**Keywords:** enumerative combinatorics, multiset, inclusion exclusion principle, constrained $k$-resolutions, cardinality of the support of the multivariate hypergeometric distribution.

## 1  Introduction

A multiset is a generalisation of a set that allows elements to appear an integer number of times, rather than being simply present or absent [2]. The number of times that an element appears in a multiset is termed its multiplicity. One multiset may be considered a sub-multiset of a second when it does not contain any element with a greater multiplicity. The cardinality of a multiset is the sum of its elements' multiplicities, the number of its elements is a distinct quantity called the dimension. *In this report, we provide a formula for the number of sub-multisets of a given multiset that have a specified cardinality.*

A concrete example of a multiset is found in a packet of candies distinguishable only by their colour. In this case, each colour of candy corresponds to a distinct element of the multiset whose multiplicity is simply the number of candies of that colour that are present. To ask how many distinct handfuls of a given size may be obtained from the packet is to ask how many multisets of a given cardinality exist that are sub-multisets of the multiset represented by the packet.

---

[*]s.ferraris@ucl.ac.uk
[†]a.mendelson@ucl.ac.uk
[‡]g.ballesio@gmail.com
[§]t.vercauteren@ucl.ac.uk

To our surprise, even after extensive searching, we were not able to find an answer to this question in the literature. While the number may be found through exhaustive enumeration of the type described in [4], this will quickly become impractical as the dimensions and multiplicities of the multisets increase. In this report, we use a proof based on the $k$-resolutions of $n$ and the inclusion-exclusion principle to derive an expression of this quantity as the sum of binomial coefficients presented in equation (10). A concise proof is provided in section 5, while a longer, step-by-step derivation is provided in sections 2 to 4.

Python code to verify the provided formula is available online at
https://github.com/SebastianoF/counting_sub_multisets.git.

## 2 Sampling from a multiset

We begin with a simple multiset that possesses 2 elements with an equal multiplicity of 5:

$$S = \{b, b, b, b, b, w, w, w, w, w\}.$$

In our candy analogy, this corresponds to a packet containing 5 black candies and 5 white ones. In this case, if we want to know how many handfuls of 5 may be produced[1], we can simply enumerate them, as has been done in the table below. Here, each row represents a handful, and the black and white circles represent the candies of the corresponding colors.

| 1 | ● | ● | ● | ● | ● |
|---|---|---|---|---|---|
| 2 | ● | ● | ● | ● | ○ |
| 3 | ● | ● | ● | ○ | ○ |
| 4 | ● | ● | ○ | ○ | ○ |
| 5 | ● | ○ | ○ | ○ | ○ |
| 6 | ○ | ○ | ○ | ○ | ○ |

The number of black candies must be an integer between 0 and 5, this also specifies the number of white candies, and therefore a total of 6 handfuls are possible. Generalising this, let the multiplicities of the elements (black and white candies) be denoted $a_1$ and $a_2$ respectively, the cardinality of the multiset (packet) be denoted $N = a_1 + a_2$, and the cardinality of the sub-multisets (handfuls) be denoted $n$. Providing $n$ is not greater than $a_1$ or $a_2$, the multiplicity of the first element in the sub-multiset may range between 0 and $n$, so there will be $n + 1$ possible outcomes.

If we allow the size of the handful to exceed the number of candies of a given colour, there is an additional constraint we must consider: it is now possible to "run out of" one colour of candy as the handful is drawn. If, instead of the original packet, we drew a handful of 5 from a smaller packet containing only 3 black candies and 4 white ones, we would no longer be able to produce the handfuls numbered 1, 2 and 6 in the table above; the number of black candies we draw must be between $n - a_2 = 1$ and $a_1 = 3$. Generalising again, a multiset with two elements of multiplicities $a_1$ and $a_2$ has the following number of sub-multisets with cardinality $n$:

$$r = \min(n, a_1) - \max(1, n - a_2) + 2.$$

While it has been straightforward to solve this problem in the case of two elements, when we increase the number of colours to a number beyond two, this can no longer be done so

---

[1] by drawing without replacement, as is usually the case with candies

intuitively. For example, we might ask how many distinct handfuls of 12 may be drawn from a packet containing 5 blue candies, 9 green ones, and 14 red ones.[2] To answer this type of question, we must reformulate the problem in terms of integer resolutions.

# 3    Reformulating the problem with $k$-resolutions

We consider the multiset $S$ of cardinality $N$ and dimension $k$ where each element $s_j$ has a multiplicity of $a_j$:

$$S = \{\underbrace{s_1, s_1, ..., s_1}_{a_1\text{-instances}}, \underbrace{s_2, s_2, ..., s_2}_{a_2\text{-instances}}, ..., \underbrace{s_k, s_k, ..., s_k}_{a_k\text{-instances}}\},$$

where

$$a_1 + a_2 + \cdots + a_k = N.$$

Providing some ordering is available for the elements, the sequence of their multiplicities $(a_1, a_2, ..., a_k)$ is all that is required to uniquely specify $S$. Any sub-multiset of $S$ with cardinality $n$ $(n \leq N)$ can be similarly specified by a sequence $(x_1, x_2, ..., x_k)$ which must satisfy the constraints

$$x_1 + x_2 + \cdots + x_k = n, \qquad 0 \leq x_j \leq a_j. \tag{1}$$

Sequences $(x_1, x_2, ..., x_k)$ that satisfy only $x_1 + x_2 + \cdots + x_k = n$ are known as $k$-resolutions of $n$, and there are established techniques to count them (see [3] or [5], for example). The strategy we use to count sequences that satisfy (1) is based on their correspondence to the $k$-resolutions of $n$ that satisfy the additional constraints $0 \leq x_j \leq a_j$.

### $k$-resolutions of $n$

We denote the number of the $k$-resolutions of $n$

$$R_k^n = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \ x_i \geq 0\}|. \tag{2}$$

For any choice of $k$ and $n$, $R_k^n$ can be found by considering a row of $n$ indistinguishable candies and $k-1$ dividers that can be placed to split the row of candies into $k$ subsets. Here, the number of candies in the $j^{th}$ subset corresponds to the integer value of the summand $x_j$. For example, the 3-resolution of 8 given by $3 + 1 + 4 = 8$ is represented

$$\bullet \quad \bullet \quad \bullet \mid \bullet \mid \bullet \quad \bullet \quad \bullet \quad \bullet$$

With this representation, the number of permutations of the full set of candies and dividers is $(n + k - 1)!$. Because we consider all candies and dividers indistinguishable, this number will overestimate the number of distinguishable arrangements by a factor of $n! \times (k - 1)!$. Correcting for this, we obtain

$$R_k^n = \frac{(n + k - 1)!}{n!(k - 1)!} = \binom{n + k - 1}{k - 1} = \binom{n + k - 1}{n}. \tag{3}$$

This is the number of sub-multisets of cardinality $n$ that a multiset has when $n$ is smaller than all $a_j$.

---

[2] If we were to draw handfuls at random, their probabilities would be described by the multivariate hypergeometric distribution [1]. The number of possible handfuls corresponds to the cardinality of its support.

## $k$-resolutions of $n$ with lower constraints

Before considering the constraints of the form $0 \le x_j \le a_j$ that occur in our original problem, it is instructive to consider another type of constraints: those of the form $x_j \ge a_j \ge 0$. We shall call these *lower constraints*. The number of $k$-resolutions that satisfy the set of lower constraints specified by the sequence $(a_1, a_2, \ldots, a_k)$ is given

$$R^n_{(a_1, a_2, \ldots, a_k)^\downarrow} = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \ x_i \ge a_i\}|.$$

To determine this quantity, we perform the substitution $y_j = x_j - a_j$ for each of the $x_j$:

$$x_1 + x_2 + \cdots + x_k = n \qquad\qquad x_i \ge a_i$$

$$x_1 - a_1 + x_2 - a_2 + \cdots + x_k - a_k = n - \sum_{j=1}^{k} a_j \qquad\qquad x_i - a_i \ge 0 \,.$$

This produces the equivalent equation

$$y_1 + y_2 + \cdots + y_k = n - \sum_{j=1}^{k} a_j \qquad\qquad y_i \ge 0 \qquad y_i = x_i - a_i \,. \qquad (4)$$

Noting that $n - \sum_{j=1}^{k} a_j$ is not negative, we can see that this equation has the same form as equation (2). By equation (3), the number of solutions to equation (4) is $R_k^{n - \sum_{j=1}^{k} a_j}$, and therefore

$$R^n_{(a_1, a_2, \ldots, a_k)^\downarrow} = \binom{n - \sum_{j=1}^{k} a_j + k - 1}{k - 1} = \binom{n - \sum_{j=1}^{k} a_j + k - 1}{n - \sum_{j=1}^{k} a_j}. \qquad (5)$$

## $k$-resolutions with upper constraints

We now consider the case that corresponds to our original problem. Here, each of the $x_j$ satisfies a constraint of the form $x_j \le a_j$. We denote the number of *upper constrained $k$-resolutions of $n$*

$$R^n_{(a_1, a_2, \ldots, a_k)^\uparrow} = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, 0 \le x_i \le a_i\}|. \qquad (6)$$

The previous strategy can not be adapted so easily to this new case. We could make the substitution

$$x_1 + x_2 + \cdots + x_k = n \qquad\qquad x_i \le a_i$$
$$x_1 + x_2 + \cdots + x_k = n \qquad\qquad a_i - x_i \ge 0$$
$$a_1 - x_1 + a_2 - x_2 + \cdots + a_k - x_k = \sum_{j=1}^{k} a_j - n \qquad\qquad a_i - x_i \ge 0$$

to obtain the form

$$y_1 + y_2 + \cdots + y_k = \sum_{j=1}^{k} a_j - n \qquad\qquad y_i = a_i - x_i \,. \qquad (7)$$

4

It might seem like we could apply the same reasoning we used for equation (4) to produce the result

$$R^n_{(a_1,a_2,\dots,a_k)\uparrow} = \binom{\sum_{j=1}^k a_j - n + k - 1}{k - 1} = \binom{\sum_{j=1}^k a_j - n + k - 1}{\sum_{j=1}^k a_j - n}, \qquad (8)$$

but *this formula is not correct.* Unlike those of equation (4), the corresponding $y_j$ of equation (7) may take negative values.

This can be seen in the counter example that follows. If we consider the problem specified by $n = 5$, $k = 3$ and $(a_1, a_2, a_3) = (2, 3, 3)$, then the possible resolutions are the 9 that follow in the form $[x_1, x_2, x_3]$:

$$[0, 2, 3], [0, 3, 2], [1, 1, 3],$$
$$[1, 2, 2], [1, 3, 1], [2, 0, 3],$$
$$[2, 1, 2], [2, 2, 1], [2, 3, 0].$$

As can be seen below, equation (8) would suggest that there are 10:

$$\binom{8 - 5 + 3 - 1}{3 - 1} = 10.$$

This has happened because we have also counted the invalid solution

$$[-1, 3, 3].$$

To find $R^n_{(a_1,a_2,\dots,a_k)\uparrow}$, the number of upper constrained $k$-resolutions of $n$, we must use a different strategy.

# 4  Solution from the inclusion-exclusion principle

We begin again with the equations specifying a $k$-resolution of $n$ as a sequence of integers $x_j$ that satisfy the original constraints:

$$x_1 + x_2 + \cdots + x_k = n, \qquad 0 \le x_i \le a_i.$$

We consider the following:

$$A = \{(x_1, x_2, \dots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \forall i \ 0 \le x_i \le a_i\}$$
$$B = \{(x_1, x_2, \dots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \exists i \ x_i \ge a_i + 1\}$$
$$C = \{(x_1, x_2, \dots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \forall i \ x_i \ge 0\}.$$

Here, $C$ is the full set of $k$-resolutions of $n$, $A$ is the set of the upper constrained resolutions (whose cardinality we wish to determine), and $B$ is the set where at least one of the $x_j$ exceeds the shared bounds $a_j$. We can see that $A = C \setminus B$ and $|A| = |C| - |B|$.

By (3), the cardinality of $C$ is

$$|C| = \binom{n + k - 1}{k - 1}.$$

To determine $|A|$, we now only need to determine $|B|$. $B$ may be written as the union

$$B = B_1 \cup B_2 \cup \cdots \cup B_k\,,$$

where

$$B_j = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, x_j \geq a_j + 1\}\,.$$

By the inclusion-exclusion principle,

$$|B_1 \cup B_2 \cup \cdots \cup B_k| = \sum_{m=1}^{k} (-1)^{m+1} \sum_{1 \leq i_1 < i_2 < \cdots < i_m \leq k} |B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m}|\,.$$

At this point, we note that the sets $B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m}$ are sets of *lower constrained* $k$-resolutions of the form

$$B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m} = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \, x_{i_l} \geq a_{i_l} + 1\,, l = 1, \ldots, m\}\,.$$

By equation (5), the value of each summand must therefore be

$$|B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m}| = \binom{n - (a_{i_1} + 1 + a_{i_2} + 1 + \cdots + a_{i_m} + 1) + k - 1}{k-1}$$
$$= \binom{n - (\sum_{l=1}^{m} a_{i_l} + m) + k - 1}{k-1}\,.$$

With this last formula we can finally determine the cardinality of $A$ (i.e., the number of upper constrained k-resolutions):

$$R^n_{(a_1, a_2, \ldots, a_k)\uparrow} = |A| = |C| - |B|$$
$$= \binom{n+k-1}{k-1} - \sum_{m=1}^{k}(-1)^{m+1} \sum_{1 \leq i_1 < i_2 < \cdots < i_m \leq k} \binom{n - \sum_{l=1}^{m} a_{i_l} - m + k - 1}{k-1}\,.$$

This can be expressed more concisely as

$$R^n_{(a_1, a_2, \ldots, a_k)\uparrow} = \sum_{L \in \mathcal{P}(I_k)} (-1)^{|L|} \binom{n + k - 1 - |L| - \sum_{i \in L} a_i}{k-1}\,, \tag{9}$$

where $\mathcal{P}(I_k)$ denotes the power set of $I_k$, and binomial coefficients of the form $\binom{\alpha}{\beta}$ with $\alpha < 0$ or $\alpha < \beta$ are taken to be zero in line with their combinatorial meaning.

# 5 Concise proof

This section provides a summary of the definitions and proofs presented in this report.

**Definition 1** *Given a finite set $A$ of cardinality $k$, a* multiset *$S$ is a collection of elements of $A$ where each element of $A$ can appear zero or more times, and where the order of the elements and appearances does not matter. The number of occurrences of an element $a \in A$ in $S$ is called its multiplicity. The sum of all a multiset's elements' multiplicities is its cardinality. The dimension of the multiset $S$ is the number of distinct elements that it contains one or more times.*

**Definition 2** *A sub-multiset $S'$ of a multiset $S$ is a multiset whose elements are all also elements of $S$. The elements' multiplicities in $S'$ must be less than or equal to their corresponding multiplicities in $S$.*

We wish to determine the number of multisets of cardinality $n$ that may be considered submulitisets of a multiset $S$ of cardinality $N$ and dimension $k$. The $k$ distinct elements have multiplicities in $S$ specified by the sequence $(a_1, a_2, \ldots, a_k)$. This problem can be formulated in terms of constrained resolutions of integers.

**Definition 3** *A $k$-resolution of $n$ is an element of the set*

$$\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \ x_i \geq 0\}.$$

*The cardinality of this set is denoted*

$$R_k^n = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \ x_i \geq 0\}|.$$

*The cardinality of the set of $k$-resolutions with lower constraints specified by a sequence $(a_1, a_2, \ldots, a_k)$ is*

$$R_{(a_1, a_2, \ldots, a_k)^\downarrow}^n = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, 0 \geq x_i \geq a_i\}|,$$

*while the cardinality of the set of $k$-resolutions with upper constraints specified by such a sequence is*

$$R_{(a_1, a_2, \ldots, a_k)^\uparrow}^n = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, x_i \leq a_i\}|.$$

The number of sub-multisets is thus equivalent to the number of $k$-resolutions of $n$ with the upper constraints $(a_1, a_2, \ldots, a_k)$.

**Lemma 1** *The cardinality of the set of $k$-resolutions with lower constraints specified by $(a_1, a_2, \cdots, a_k)$ is given*

$$R_{(a_1, a_2, \ldots, a_k)^\downarrow}^n = \binom{n - \sum_{j=1}^k a_j + k - 1}{k - 1}.$$

**Proof:**

$$R_{(a_1, a_2, \ldots, a_k)^\downarrow}^n = |\{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, x_i \geq a_i\}|$$

$$= |\{(x_1, x_2, \ldots, x_k) \mid x_1 - a_1 + x_2 - a_2 + \cdots + x_k - a_k = n - \sum_{j=1}^k a_j, x_i \geq a_i\}|$$

$$= |\{(x_1, x_2, \ldots, x_k) \mid y_1 + y_2 + \cdots + y_k = n - \sum_{j=1}^k a_j, \ y_i \geq 0\}|$$

and so using the standard result for unconstrained $k$-resolutions,

$$= \binom{n - \sum_{j=1}^k a_j + k - 1}{k - 1}.$$

$\square$

**Theorem 1** *Indicating $\{1, 2, \cdots, k\}$ with $I_k$, the cardinality of the set of upper constrained k-resolutions is*

$$R^n_{(a_1, a_2, \ldots, a_k)\uparrow} = \sum_{L \in \mathcal{P}(I_k)} (-1)^{|L|} \binom{n + k - 1 - |L| - \sum_{i \in L} a_i}{k - 1}, \tag{10}$$

*where $\mathcal{P}(I_k)$ is the power set of $I_k$.*

**Proof:** We define the sets $A$, $B$ and $C$ as follows:

$$A = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \forall i \ \ 0 \le x_i \le a_i\}$$
$$B = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \exists i \ \ x_i \ge a_i + 1\}$$
$$C = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, \text{ such that } \forall i \ \ x_i \ge 0\}.$$

From this, it follows that

$$A \cup B = C \qquad A \cap B = \emptyset \qquad |A| = |C| - |B| \qquad |C| = \binom{n + k - 1}{k - 1}.$$

The set $B$ can be written as the union of subsets

$$B = B_1 \cup B_2 \cup \cdots \cup B_k,$$

where

$$B_j = \{(x_1, x_2, \ldots, x_k) \mid x_1 + x_2 + \cdots + x_k = n, x_j \ge a_j + 1\}.$$

Due to the inclusion-exclusion principle,

$$|B_1 \cup B_2 \cup \cdots \cup B_k| = \sum_{m=1}^{k} (-1)^{m+1} \sum_{1 \le i_1 < i_2 < \cdots < i_j \le k} |B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m}|.$$

From the previous lemma, it follows that

$$|B_{i_1} \cap B_{i_2} \cap \cdots \cap B_{i_m}| = \binom{n - (a_{i_1} + 1 + a_{i_2} + 1 + \cdots + a_{i_m} + 1) + k - 1}{k - 1},$$

and therefore

$$R^n_{(a_1, a_2, \ldots, a_k)\uparrow} = \binom{n + k - 1}{k - 1} - \sum_{m=1}^{k} (-1)^{m+1} \sum_{1 \le i_1 < i_2 < \cdots < i_j \le k} \binom{n - \sum_{l=1}^{m} a_{i_l} - m + k - 1}{k - 1}.$$

$\square$

# Acknowledgment

# References

[1] Berkopec, Aleš. HyperQuick algorithm for discrete hypergeometric distribution. Journal of Discrete Algorithms 5.2 (2007): 341-347.

[2] Bona, Miklos. Introduction to enumerative combinatorics. McGraw-Hill Higher Education, 2007.

[3] C. Mariconda, A. Tonolo, *Calcolo Discreto*, Apogeo 2012.

[4] Hage, Jurriaan. Enumerating submultisets of multisets. Information processing letters 85.4 (2003): 221-226.

[5] Patashnik, RL Graham DE Knuth O. Concrete Mathematics-A Foundation for Computer Science. (1989).

[6] Ross, Sheldon M. Introduction to probability and statistics for engineers and scientists. Academic Press, 2014.