



SCHOOL OF MATHEMATICS AND STATISTICS

LEVEL-5 MSCI PROJECT

Size-Biased Sampling from a Population of Clusters

Author:

James Marriner
2554917M

Supervisors:

Dr Alexey Lindo
Dr Serik Sagitov

Declaration of Originality

I confirm that this assignment is my own work and that I have:

- Read and understood the guidance on plagiarism in the Student Handbook, including the University of Glasgow Statement on Plagiarism
- Not made use of the work of any other student(s) past or present without acknowledgment. This includes any of my own work, that has been previously, or concurrently, submitted for assessment, either at this or any other educational institution, including school
- Not sought or used the services of any professional agencies to produce this work
- In addition, I understand that any false claim in respect of this work will result in disciplinary action in accordance with University regulations
- I am aware of and understand the University's policy on plagiarism and I certify that this assignment is my own work, except where indicated by referencing, and that I have followed the good academic practices noted above.

April 6, 2025

Abstract

We aim to understand how the number of classes of a given size change as samples are drawn from a population partitioned into distinct classes. Derivations of the expectations, variances and covariances are given as well as an example to illustrate their form. We show how this problem appears within the common species richness estimation problem and thus connect our results. We then derive asymptotics results for both classes of a given size, and the number of classes when sampling without replacement.

Contents

1	Introduction	2
2	The Dilution Process	3
2.1	PCR	3
2.2	A Clustered Population	3
3	Dilution Process Dynamics	4
3.1	The Simplest Example	5
3.2	Indicator Random Variables	6
3.3	Second Moments	7
3.4	Variance	9
3.5	Covariances	9
3.6	An Illustrative Example	11
3.7	Population Results	14
3.8	Average Cluster Size	15
3.8.1	Simulation Results	16
4	Sampling with Replacement	17
4.0.1	Derivation of Higher Moments	17
5	Species Richness Estimation	19
6	Asymptotic Results	20
6.1	Extra Notation	22
6.2	Limit Theorems	22
7	Conclusion	27
8	Appendix	28
8.1	Recurrence Relations	28
8.2	Population Variance and Covariance	31
8.3	Variance of Species Estimators	32
8.4	Harmonic Inequality	33

1 Introduction

The problem of estimating the number of distinct classes of objects that partition a population has been considered in statistical settings as early as Fisher et al. in 1943 [1], who sought a connection between the number of species and the number of specimens using a sample of butterflies. This question has wide ranging applications from ecologists estimating species richness [2], which is a strong measure of a region's biodiversity, to linguists who have tried to estimate Shakespeare's vocabulary size given the words used in his work [3].

Generally in these contexts, the goal is to estimate the number of classes using a sample drawn with or without replacement (Poisson frequency models are also used). The number of classes seen in the sample is generally not a good estimator in this case, as the sample provides no information about the number of species left unobserved in the population. Resulting estimates do improve for larger samples, however, the problem can be particularly hard, especially in the case when there exist many small and hence rare, classes.

For this problem, numerous methods and estimators have been proposed, surveys and comparisons of which can be found in [4] or [5]. Each of these approaches are uniquely tuned towards class estimation, but distinctly less focus has been placed on estimating the number of classes of a given size.

$$\hat{C}_{\text{GOODMAN}} = c + \sum_{j=1}^n (-1)^{j+1} \frac{\binom{N-n+j-1}{j}}{\binom{n}{j}} c_j \quad \hat{C}_{\text{SHLOSSER}} = c + \frac{c_1}{\sum_{i=1}^n iq(1-q)^{n-1}} \cdot \sum_{j=1}^n (1-q)^j c_j$$

Figure 1: Goodman's and Shlosser's Estimators for the Number of Classes given a Sample of Size n . c is the number of observed classes, c_j the number of observed classes of size j and $q = \frac{n}{N}$

In the context of species richness estimation, estimating the number of species of a particular size has little application, as this represents an arbitrary grouping of animals, however, there is some existing usage for this estimate as seen in [6], when estimating words of a particular size in an unseen dictionary.

By examining the size-biased sampling at the heart of the dilution process within PCR applied to a cancer biopsy, we uncover a particular focus on estimating the number of groups with the same size, and therefore, aim to estimate these values as increasing numbers of samples are taken. In addition, in the case of sampling without replacement, we investigate the relationship between our estimates and the remaining population from which the samples are drawn.

Finally, we connect our estimates to the familiar problem of estimating the total number of classes to which the existing body of literature is devoted. Having recovered some standard results of Hurlbert [7], we seek to derive asymptotics for both our group estimates in the dilution process, and overall estimates in the general case.

2 The Dilution Process

2.1 PCR

Polymerase Chain Reaction, henceforth PCR, is a method used to analyse DNA in a variety of settings, such as from the small sample of a cancer biopsy. To improve the conclusive detection of abnormalities such as cancer, as opposed to typical random mutations, part of PCR involves multiple steps of amplifying and diluting the samples before final inference is performed. In the PCR setting we are interested in, the process is as follows:

1. An initial sample of N molecules is taken,
2. The 1st amplification raises N to $\approx 100,000$,
3. The 1st dilution decreases N to $\approx 80,000$,
4. The 2nd amplification raises N to $\approx 100,000,000,000$,
5. The 2nd dilution decreases N to $\approx 100,000$. The final analysis is performed on this set of molecules.

The effect of the dilution process from a statistical perspective is poorly understood, especially given the unique characteristics that the molecules possess. Our goal is to explore this action from the perspective of a stochastic process.

2.2 A Clustered Population

Suppose that a population of size N , such as from a cancer biopsy, is composed of molecules of different types, denoted by q for $q = 1, \dots, Q$, where Q is the number of different types. Additionally, suppose that molecules of the same type belong to clusters, or groups, of size $c = 1, \dots, M_q$ within this population, where M_q is the maximum cluster size for molecules of type q .

Such a population may be described by the proportion of molecules which belong to clusters of a particular size and type. Let;

- $\alpha_{q,c}$ be the proportion of molecules of type q in clusters of size c ,
- $N_{q,c}$ be the number of molecules of type q belonging to clusters of size c .

Where $0 \leq \alpha_{q,c} \leq 1$ for all q, c and

$$\sum_{q=1}^Q \sum_{c=1}^{M_q} \alpha_{q,c} = 1$$

We also have that;

$$N_{q,c} = \alpha_{q,c} \cdot N \quad \text{and} \quad \sum_{q=1}^Q \sum_{c=1}^{M_q} N_{q,c} = N$$

We likewise define $\alpha_{\bullet,c} = \sum_{q=1}^Q \alpha_{q,c}$ as the overall proportion of molecules belonging to clusters of size c and $\alpha_{q,\bullet} = \sum_{c=1}^{M_q} \alpha_{q,c}$ as the proportion of molecules of type q . The corresponding numbers of molecules in these categories, $N_{\bullet,c}$ and $N_{q,\bullet}$, are then defined accordingly.

We are interested in the effect of the dilution process where sampling from such a population is conducted without replacement. Primarily, this is interesting as the distribution of molecules in clusters

of various sizes will change as more samples are drawn.

During each dilution phase, we take n samples from our ‘population’ or equivalently sample a single molecule at each discrete time point $t = 1, \dots, n$. If we sample from a cluster of size $c \geq 2$, then the remaining molecules will reduce to a cluster of size $c - 1$ within the population. Likewise, we assume that molecules are uniquely identifiable by the cluster they originally belonged to. As a result, if we sample more than once from the same original cluster, the molecules will be regrouped in the sample.

Samples Taken	Population	Sample
0		
1		
2		
3		

Table 1: Repeatedly sampling from a cluster of size 4

We primarily aim to understand how the distribution of cluster sizes in the sample evolves as samples are taken, but the distribution within the remaining population is also a secondary concern.

Functionally, we will see that the molecule type q has little effect on the dynamics of the process as molecules of distinct types never appear in the same clusters. Furthermore, we can consider a two-step sampling process of:

1. Choosing a molecule type q with probability equal to the proportion of type q molecules remaining in the population,
2. Sample a molecule without replacement from the grouped molecules of type q .

In this case, step 2 corresponds exactly to the dilution process with $Q = 1$ and thus, multiple molecule types should only proportionally slow the behaviour of the process as samples are taken. Nonetheless, as the molecules within PCR exhibit different biological types, we include this feature.

3 Dilution Process Dynamics

For the sampling procedure described above, let;

1. $X_{q,c}(n)$ be the number of clusters containing c molecules of type q in the population after n samples are taken,
2. $S_{q,c}(n)$ be the number of clusters containing c molecules of type q in the sample after n molecules are drawn without replacement.

Trivially we then note that for all q and c , the number of clusters in the original population is given by;

$$S_{q,c}(N) = X_{q,c}(0) = \frac{N_{q,c}}{c}$$

The equality holds as if all N molecules are sampled, all original clusters are reformed, and the final sample resembles exactly the initial population.

3.1 The Simplest Example

To illustrate the process more clearly, we consider the simplest case where $Q = 1$ and $M = 2$. If we consider sequential samples being taken, we deduce the following relations for the population when sampling the $(n + 1)$ 'th molecule;

$$X_{q=1,c=1}(n+1) = \begin{cases} X_{1,1}(n) + 1, & \text{if we sample from a cluster of size 2} \\ X_{1,1}(n) - 1, & \text{if we sample from a cluster of size 1} \end{cases}$$

$$X_{1,2}(n+1) = \begin{cases} X_{1,2}(n) - 1, & \text{if we sample from a cluster of size 2} \\ X_{1,2}(n), & \text{if we sample from a cluster of size 1} \end{cases}$$

Similarly, for the number of clusters in the sample, we have:

$$S_{1,1}(n+1) = \begin{cases} S_{1,1}(n) - 1, & \text{if the sampled molecule forms a cluster of size 2} \\ S_{1,1}(n) + 1, & \text{otherwise} \end{cases}$$

$$S_{1,2}(n+1) = \begin{cases} S_{1,2}(n) + 1, & \text{if the sampled molecule forms a cluster of size 2} \\ S_{1,2}(n), & \text{otherwise} \end{cases}$$

With these descriptions, recurrence relations can be derived for the expected numbers of clusters. When solved, they return;

$$\mathbb{E}[X_{1,1}(n)] = \frac{N-n}{N-1}(\alpha_{1,1}(N-1) + \alpha_{1,2}n)$$

$$\mathbb{E}[X_{1,2}(n)] = \frac{\alpha_{1,2}(N-n)(N-n-1)}{2(N-1)}$$

$$\mathbb{E}[S_{1,1}(n)] = \frac{n}{N-1}(\alpha_{1,1}(N-1) + \alpha_{1,2}(N-n))$$

$$\mathbb{E}[S_{1,2}(n)] = \frac{\alpha_{1,2}n(n-1)}{2(N-1)}$$

As an initial observation, we note that for $M_1 = 2$, these expressions are quadratic in n and appear to display the following symmetry for clusters of size 1 and 2:

$$\mathbb{E}[X_{1,c}(n)] = \mathbb{E}[S_{1,c}(N-n)]$$

Additionally, recalling that we can define the initial number of clusters of size c by:

$$X_{q,c}(0) = \frac{\alpha_{q,c}N}{c}$$

We see that these may be rephrased as:

$$\mathbb{E}[X_{1,1}(n)] = X_{1,1}(0) \frac{\binom{N-1}{N-n-1}}{\binom{N}{n}} + X_{1,2}(0) \frac{\binom{2}{1} \binom{N-2}{N-n-1}}{\binom{N}{n}}$$

$$\mathbb{E}[X_{1,2}(n)] = X_{1,2}(0) \frac{\binom{N-2}{N-n-1}}{\binom{N}{n}}$$

$$\mathbb{E}[S_{1,1}(n)] = X_{1,1}(0) \frac{\binom{N-1}{n-1}}{\binom{N}{n}} + X_{1,2}(0) \frac{\binom{2}{1} \binom{N-2}{n-1}}{\binom{N}{n}}$$

$$\mathbb{E}[S_{1,2}(n)] = X_{1,2}(0) \frac{\binom{N-2}{n-1}}{\binom{N}{n}}$$

These expressions resemble hypergeometric probabilities which we may expect when sampling from a varied population without replacement. Whilst these recurrence relations do allow us to find an expression for the expectation in the general case (see 8.1), when looking for higher moments, this approach is infeasible and a different strategy is needed.

3.2 Indicator Random Variables

Returning to the general case, we consider how a cluster of size c can form in the sample. These c molecules may correspond to a cluster originally of size c in the population, however, they could also correspond to a cluster of size $c + 1, c + 2, \dots$ from which only c molecules have been sampled. This illustrates that $S_{q,c}(n)$ can draw contributions from any cluster of size $c, c + 1, \dots, M_q$.

For each $k = 1, \dots, M_q$ we label the clusters by $i_k = 1, \dots, X_{q,k}(0)$. This allows us to introduce:

$$\mathbb{1}_{i_k, k, q}^c(n) = \begin{cases} 1, & \text{if the } i_k \text{'th cluster has been sampled from } c \text{ times in } n \text{ samples} \\ 0, & \text{otherwise} \end{cases}$$

where by the i_k 'th cluster, we specifically mean of original size k and type q . The numbers of clusters in the sample can then be defined using these indicators:

$$S_{q,c}(n) = \sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k, k, q}^c(n) \quad (1)$$

Therefore, by linearity of expectation;

$$\begin{aligned} \mathbb{E}[S_{q,c}(n)] &= \sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{E}[\mathbb{1}_{i_k, k, q}^c(n)] \\ &= \sum_{k=c}^{M_q} X_{q,k}(0) \mathbb{P}(\mathbb{1}_{i_k, k, q}^c(n) = 1) \end{aligned}$$

We have dispensed with summation over i_k , since this probability is the same for all clusters of size k . To determine this probability, we note that there are:

- $\binom{N}{n}$ ways to draw a sample of size n without replacement from the population,
- $\binom{k}{c}$ ways to sample the necessary c molecules from a cluster of size k ,
- $\binom{N-k}{n-c}$ ways to sample the remaining molecules which make a sample of size n .

Therefore, we find that;

$$\mathbb{P}(\mathbb{1}_{i_k, k, q}^c(n) = 1) = \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}}$$

which we recognise as a hypergeometric probability, matching our intuition about the process. Hence for all $q = 1, \dots, Q$ and $c = 1, \dots, M_q$;

$$\mathbb{E}[S_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}}$$

This can also be written compactly for a single molecule type q by defining the vector of expectations:

$$\underline{\mathbb{E}[S_q(n)]} = (\mathbb{E}[S_{q,1}(n)], \mathbb{E}[S_{q,2}(n)], \dots, \mathbb{E}[S_{q,M_q}(n)])^T$$

As well as the initial distribution vector for type q molecules:

$$\underline{X}_q(0) = (X_{q,1}(0), X_{q,2}(0), \dots, X_{q,M_q}(0))^T$$

and lastly the matrix $A(n) = (a_{i,j}(n))_{i,j=1,\dots,M_q}$ defined by:

$$a_{i,j}(n) = \begin{cases} \frac{\binom{j}{i} \binom{N-j}{n-i}}{\binom{N}{n}}, & \text{if } i \leq j \\ 0, & \text{otherwise} \end{cases}$$

With these tools, we then have that:

$$\underline{\mathbb{E}[S_q(n)]} = A(n) \underline{X}_q(0)$$

As illustrated in [6], if it is possible to invert A , which may be done with stability if $\frac{n}{N}$ is large enough, then the initial distribution $\underline{X}_q(0)$ can be estimated for a given sample.

The same indicators can be used to define the number of clusters remaining in the population after n samples. For a cluster to contribute to $X_{q,c}(n)$, it must be originally of size $k \geq c$ and have been sampled from $k - c$ times in n samples. Using the same indicators as before, we find that;

$$X_{q,c}(n) = \sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{I}_{i_k,k,q}^{k-c}(n)$$

So that;

$$\mathbb{E}[X_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-(k-c)}}{\binom{N}{n}}$$

Noting that:

$$\binom{k}{c} = \binom{k}{k-c}, \quad \binom{N-k}{n-(k-c)} = \binom{N-k}{N-n-c}, \quad \text{and} \quad \binom{N}{n} = \binom{N}{N-n}$$

we may rewrite the expectation as;

$$\mathbb{E}[X_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{N-n-c}}{\binom{N}{N-n}} = \mathbb{E}[S_{q,c}(N-n)]$$

which confirms the symmetry identified earlier. Additionally, note that q only impacts these expectations via the initial proportion $\alpha_{q,\bullet}$ and therefore needs little consideration when interpreting these results.

3.3 Second Moments

Defining the number of clusters as the sum of indicators is also convenient for computing higher order moments. For example, using (1) we have that;

$$\begin{aligned} S_{q,c}(n)^2 &= \left(\sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{I}_{i_k,k,q}^c(n) \right)^2 \\ &= \sum_{k=c}^{M_q} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{I}_{i_k,k,q}^c(n) \right)^2 + 2 \sum_{k=c}^{M_q-1} \sum_{l=k+1}^{M_q} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{I}_{i_k,k,q}^c(n) \right) \left(\sum_{i_l=1}^{X_{q,l}(0)} \mathbb{I}_{i_l,l,q}^c(n) \right) \end{aligned} \quad (2)$$

The first component of (2) is given by;

$$\begin{aligned} \sum_{k=c}^{M_q} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k, k, q}^c(n) \right)^2 &= \sum_{k=c}^{M_q} \left[\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k, k, q}^c(n)^2 + 2 \sum_{i_k=1}^{X_{q,k}(0)-1} \sum_{j_k=i_k+1}^{X_{q,c}(k)} \mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{j_k, k, q}^c(n) \right] \\ &= \sum_{k=c}^{M_q} \left[X_{q,k}(0) \mathbb{1}_{i_k, k, q}^c(n) + X_{q,k}(0)(X_{q,k}(0)-1) \mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{j_k, k, q}^c(n) \right] \end{aligned}$$

since as mentioned previously, the indicators are identical across i_k . The second component is;

$$2 \sum_{k=c}^{M_q-1} \sum_{l=k+1}^{M_q} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k, k, q}^c(n) \right) \left(\sum_{i_l=1}^{X_{q,l}(0)} \mathbb{1}_{i_l, l, q}^c(n) \right) = 2 \sum_{k=c}^{M_q-1} \sum_{l=k+1}^{M_q} (X_{q,k}(0) \mathbb{1}_{i_k, k, q}^c(n)) (X_{q,l}(0) \mathbb{1}_{i_l, l, q}^c(n))$$

Hence, to determine the second moment, we need only compute the expectation of these indicator products. Firstly, we have the product of indicators representing clusters of the same original size:

$$\mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{j_k, k, q}^c(n) = \begin{cases} 1, & \text{if the } i_k\text{'th and } j_k\text{'th cluster have each been sampled from } c \\ & \text{times in } n \text{ samples} \\ 0, & \text{otherwise} \end{cases}$$

where the i_k 'th and j_k 'th cluster are distinct and both of type q . Therefore, we deduce:

$$\mathbb{E}[\mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{j_k, k, q}^c(n)] = \frac{\binom{k}{c} \binom{k}{c} \binom{N-2k}{n-2c}}{\binom{N}{n}}$$

which is a multivariate hypergeometric probability. This is a particular case of the more general indicator pair:

$$\mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{i_l, l, q}^c(n) = \begin{cases} 1, & \text{if the } i_k\text{'th and } i_l\text{'th cluster have each been sampled from } c \\ & \text{times in } n \text{ samples} \\ 0, & \text{otherwise} \end{cases}$$

where k and l are distinct original sizes. We likewise find;

$$\mathbb{E}[\mathbb{1}_{i_k, k, q}^c(n) \mathbb{1}_{i_l, l, q}^c(n)] = \frac{\binom{k}{c} \binom{l}{c} \binom{N-k-l}{n-2c}}{\binom{N}{n}}$$

Combining these results, we find the second moment as;

$$\begin{aligned} \mathbb{E}[S_{q,c}(n)^2] &= \sum_{k=c}^{M_q} \left[X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} + X_{q,k}(0)(X_{q,k}(0)-1) \frac{\binom{k}{c} \binom{k}{c} \binom{N-2k}{n-2c}}{\binom{N}{n}} \right] \\ &\quad + 2 \sum_{k=c}^{M_q-1} \sum_{l=k+1}^{M_q} X_{q,k}(0) X_{q,l}(0) \frac{\binom{k}{c} \binom{l}{c} \binom{N-k-l}{n-2c}}{\binom{N}{n}} \end{aligned}$$

or, alternatively,

$$\mathbb{E}[S_{q,c}(n)^2] = \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} - \frac{\binom{k}{c} \binom{k}{c} \binom{N-2k}{n-2c}}{\binom{N}{n}} \right] + \sum_{k=c}^{M_q} \sum_{l=c}^{M_q} X_{q,k}(0) X_{q,l}(0) \frac{\binom{k}{c} \binom{l}{c} \binom{N-k-l}{n-2c}}{\binom{N}{n}} \quad (3)$$

3.4 Variance

To find the variance, it remains only to find the squared expectation. We have;

$$\mathbb{E}[S_{q,c}(n)]^2 = \left(\sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} \right)^2 = \sum_{k=c}^{M_q} \sum_{l=c}^{M_q} X_{q,k}(0) X_{q,l}(0) \frac{\binom{k}{c} \binom{l}{c} \binom{N-k}{n-c} \binom{N-l}{n-c}}{\binom{N}{n}^2}$$

so find the variance as;

$$\begin{aligned} \text{Var}(S_{q,c}(n)) &= \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} - \frac{\binom{k}{c} \binom{k}{c} \binom{N-2k}{n-2c}}{\binom{N}{n}} \right] \\ &\quad + \sum_{k=c}^{M_q} \sum_{l=c}^{M_q} X_{q,k}(0) X_{q,l}(0) \left[\frac{\binom{k}{c} \binom{l}{c} \binom{N-k-l}{n-2c}}{\binom{N}{n}} - \frac{\binom{k}{c} \binom{l}{c} \binom{N-k}{n-c} \binom{N-l}{n-c}}{\binom{N}{n}^2} \right] \end{aligned} \quad (4)$$

3.5 Covariances

The covariance between the number of clusters of different sizes can then also be found in the same way. Consider distinct cluster sizes c_1 and c_2 where, without loss of generality, we take $c_1 < c_2$. We start by computing the joint expectation.

$$\begin{aligned} S_{q,c_1}(n) S_{q,c_2}(n) &= \left(\sum_{k=c_1}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \right) \left(\sum_{k=c_2}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_2}(n) \right) \\ &= \sum_{k=c_2}^{M_q} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \right) \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_2}(n) \right) + \sum_{k \neq l} \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \right) \left(\sum_{i_l=1}^{X_{q,l}(0)} \mathbb{1}_{i_l,l,q}^{c_2}(n) \right) \end{aligned} \quad (5)$$

Compared to computation of the second moment, we must be more careful when evaluating the first component of the joint number of clusters. This term is:

$$\begin{aligned} &\left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \right) \left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_2}(n) \right) \\ &= \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \mathbb{1}_{i_k,k,q}^{c_2}(n) + 2 \sum_{i_k=1}^{X_{q,k}(0)-1} \sum_{j_k=i_k+1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \mathbb{1}_{j_k,k,q}^{c_2}(n) \end{aligned}$$

The reason for caution is due to the first indicator product;

$$\mathbb{1}_{i_k,k,q}^{c_1}(n) \mathbb{1}_{i_k,k,q}^{c_2}(n) = \begin{cases} 1, & \text{if the } i_k \text{'th cluster has been sampled from exactly } c_1 \text{ and } c_2 \\ & \text{times in } n \text{ samples} \\ 0, & \text{otherwise} \end{cases}$$

Sampling exactly c_1 and exactly c_2 times from the same cluster are mutually exclusive events, so this product must equal zero regardless of n .

For the second indicator product within the same cluster size k , we have;

$$\mathbb{1}_{i_k,k,q}^{c_1}(n) \mathbb{1}_{j_k,k,q}^{c_2}(n) = \begin{cases} 1, & \text{if the } i_k \text{'th cluster has been sampled from } c_1 \text{ times and the} \\ & \text{ } j_k \text{'th cluster has been sampled from } c_2 \text{ times in } n \text{ samples} \\ 0, & \text{otherwise} \end{cases}$$

Which, by considering the probability of this event, we find as;

$$\mathbb{E}[\mathbb{1}_{i_k,k,q}^{c_1}(n)\mathbb{1}_{j_k,k,q}^{c_2}(n)] = \frac{\binom{k}{c_1}\binom{k}{c_2}\binom{N-2k}{n-c_1-c_2}}{\binom{N}{n}} \quad (6)$$

Lastly, the component of (5) describing distinct original cluster sizes k and l contains the term;

$$\left(\sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{c_1}(n) \right) \left(\sum_{i_l=1}^{X_{q,l}(0)} \mathbb{1}_{i_l,l,q}^{c_2}(n) \right) = X_{q,k}(0)X_{q,l}(0)\mathbb{1}_{i_k,k,q}^{c_1}(n)\mathbb{1}_{i_l,l,q}^{c_2}(n)$$

which is a slightly more general version of (6) so that;

$$\mathbb{E}[\mathbb{1}_{i_k,k,q}^{c_1}(n)\mathbb{1}_{i_l,l,q}^{c_2}(n)] = \frac{\binom{k}{c_1}\binom{l}{c_2}\binom{N-k-l}{n-c_1-c_2}}{\binom{N}{n}}$$

Combining these results, we find that the joint expectation is given by;

$$\mathbb{E}[S_{q,c_1}(n)S_{q,c_2}(n)] = \sum_{k=c_2}^{M_q} X_{q,k}(0)(X_{q,k}(0) - 1) \frac{\binom{k}{c_1}\binom{k}{c_2}\binom{N-2k}{n-c_1-c_2}}{\binom{N}{n}} + \sum_{k \neq l} X_{q,k}(0)X_{q,l}(0) \frac{\binom{k}{c_1}\binom{l}{c_2}\binom{N-k-l}{n-c_1-c_2}}{\binom{N}{n}}$$

or, alternatively,

$$\mathbb{E}[S_{q,c_1}(n)S_{q,c_2}(n)] = \sum_{k=c_1}^{M_q} \sum_{l=c_2}^{M_q} X_{q,k}(0)X_{q,l}(0) \frac{\binom{k}{c_1}\binom{l}{c_2}\binom{N-k-l}{n-c_1-c_2}}{\binom{N}{n}} - \sum_{k=c_2}^{M_q} X_{q,k}(0) \frac{\binom{k}{c_1}\binom{k}{c_2}\binom{N-2k}{n-c_1-c_2}}{\binom{N}{n}} \quad (7)$$

Lastly, to compute the covariance, we need the product of expectations. Again assuming $c_1 < c_2$, this is given by;

$$\begin{aligned} \mathbb{E}[S_{q,c_1}(n)]\mathbb{E}[S_{q,c_2}(n)] &= \left(\sum_{k=c_1}^{M_q} X_{q,k}(0) \frac{\binom{k}{c_1}\binom{N-k}{n-c_1}}{\binom{N}{n}} \right) \left(\sum_{l=c_2}^{M_q} X_{q,l}(0) \frac{\binom{l}{c_2}\binom{N-l}{n-c_2}}{\binom{N}{n}} \right) \\ &= \sum_{k=c_1}^{M_q} \sum_{l=c_2}^{M_q} X_{q,k}(0)X_{q,l}(0) \frac{\binom{k}{c_1}\binom{N-k}{n-c_1}\binom{l}{c_2}\binom{N-l}{n-c_2}}{\binom{N}{n}^2} \end{aligned}$$

and thus, we find the covariance as:

$$\begin{aligned} \text{Cov}(S_{q,c_1}(n), S_{q,c_2}(n)) &= \sum_{k=c_1}^{M_q} \sum_{l=c_2}^{M_q} X_{q,k}(0)X_{q,l}(0) \left[\frac{\binom{k}{c_1}\binom{l}{c_2}\binom{N-k-l}{n-c_1-c_2}}{\binom{N}{n}} - \frac{\binom{k}{c_1}\binom{N-k}{n-c_1}\binom{l}{c_2}\binom{N-l}{n-c_2}}{\binom{N}{n}^2} \right] \\ &\quad - \sum_{k=c_2}^{M_q} X_{q,k}(0) \frac{\binom{k}{c_1}\binom{k}{c_2}\binom{N-2k}{n-c_1-c_2}}{\binom{N}{n}} \end{aligned}$$

We note that this negative effect occurs due to discussing the same specific cluster of the original population. We infer then that the result should not be the same when dealing with distinct molecule types q_1 and q_2 .

Assuming still that, $c_1 < c_2$ but we have distinct molecule types q_1, q_2 , we find the joint cluster number as;

$$S_{q_1,c_1}(n)S_{q_2,c_2}(n) = \left(\sum_{k=c_1}^{M_{q_1}} \sum_{i_k=1}^{X_{q_1,k}(0)} \mathbb{1}_{i_k,q_1,k}^{c_1}(n) \right) \left(\sum_{j=c_2}^{M_{q_2}} \sum_{i_j=1}^{X_{q_2,j}(0)} \mathbb{1}_{i_j,q_2,j}^{c_2}(n) \right)$$

as we are in no risk of describing the same cluster, we can dispense with summation over i_k .

$$S_{q_1,c_1}(n)S_{q_2,c_2}(n) = \left(\sum_{k=c_1}^{M_{q_1}} X_{q_1,k}(0)\mathbb{I}_{i_k,q_1,k}^c(n) \right) \left(\sum_{k=c_2}^{M_{q_2}} X_{q_2,k}(0)\mathbb{I}_{i_j,q_2,j}^{c_2}(n) \right)$$

and hence find the joint expectation as;

$$\mathbb{E}[S_{q_1,c_1}(n)S_{q_2,c_2}(n)] = \sum_{k=c_1}^{M_{q_1}} \sum_{l=c_2}^{M_{q_2}} X_{q_1,k}(0)X_{q_2,l}(0) \frac{\binom{k}{c_1} \binom{l}{c_2} \binom{N-k-l}{n-c-l}}{\binom{N}{n}} \quad (8)$$

With the only difference between (7) and (8) being the negative sum, we conclude that the covariance for any two molecule types q_1, q_2 and cluster sizes $c_1 < c_2$ is:

$$\begin{aligned} \text{Cov}(S_{q_1,c_1}(n), S_{q_2,c_2}(n)) &= \sum_{k=c_1}^{M_{q_1}} \sum_{l=c_2}^{M_{q_2}} X_{q_1,k}(0)X_{q_2,l}(0) \left[\frac{\binom{k}{c_1} \binom{l}{c_2} \binom{N-k-l}{n-c_1-c_2}}{\binom{N}{n}} - \frac{\binom{k}{c_1} \binom{N-k}{n-c_1} \binom{l}{c_2} \binom{N-l}{n-c_2}}{\binom{N}{n}^2} \right] \\ &\quad - \delta_{q_1,q_2} \sum_{k=c_2}^{M_{q_1}} X_{q_1,k}(0) \frac{\binom{k}{c_1} \binom{k}{c_2} \binom{N-2k}{n-c_1-c_2}}{\binom{N}{n}} \end{aligned}$$

where:

$$\delta_{q_1,q_2} = \begin{cases} 1, & \text{if } q_1 = q_2 \\ 0, & \text{otherwise} \end{cases}$$

is the standard Kronecker delta function.

3.6 An Illustrative Example

To illustrate these results, we conduct a simulation with the following parameters for $n = 1, \dots, N$:

$$Q = 1, \quad N = 100,000, \quad M_1 = 8$$

We define the distribution of clusters as roughly decreasing so that fewer molecules are contained in the largest clusters.

$$\alpha_{1,1} = 0.23, \quad \alpha_{1,2} = 0.19, \quad \alpha_{1,3} = 0.15, \quad \alpha_{1,4} = 0.11, \quad \alpha_{1,5} = 0.14, \quad \alpha_{1,6} = 0.10, \quad \alpha_{1,7} = 0.06, \quad \alpha_{1,8} = 0.02 \quad (9)$$

Beginning with the expected numbers of clusters in the sample, we have:

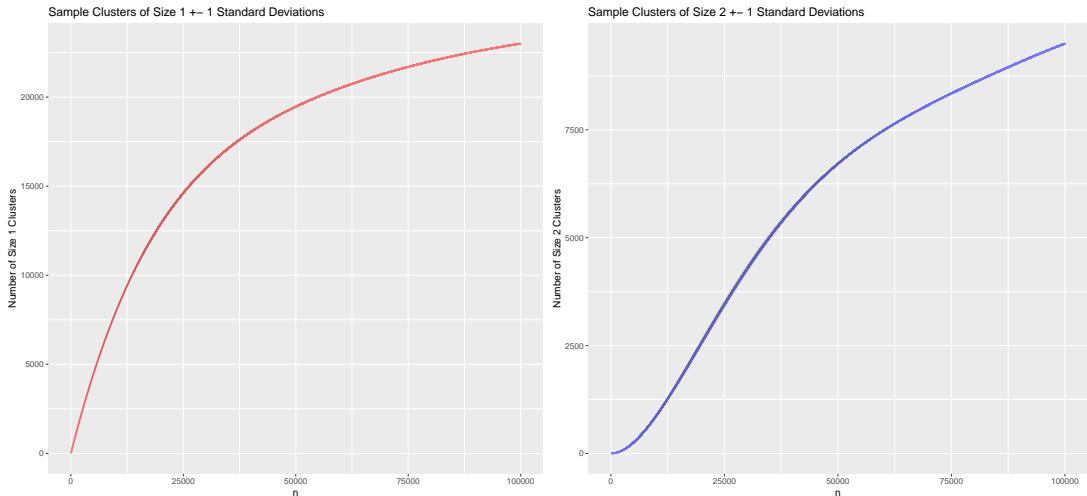


Figure 2: Expectations ± 1 Standard Deviation for $c = 1, 2$

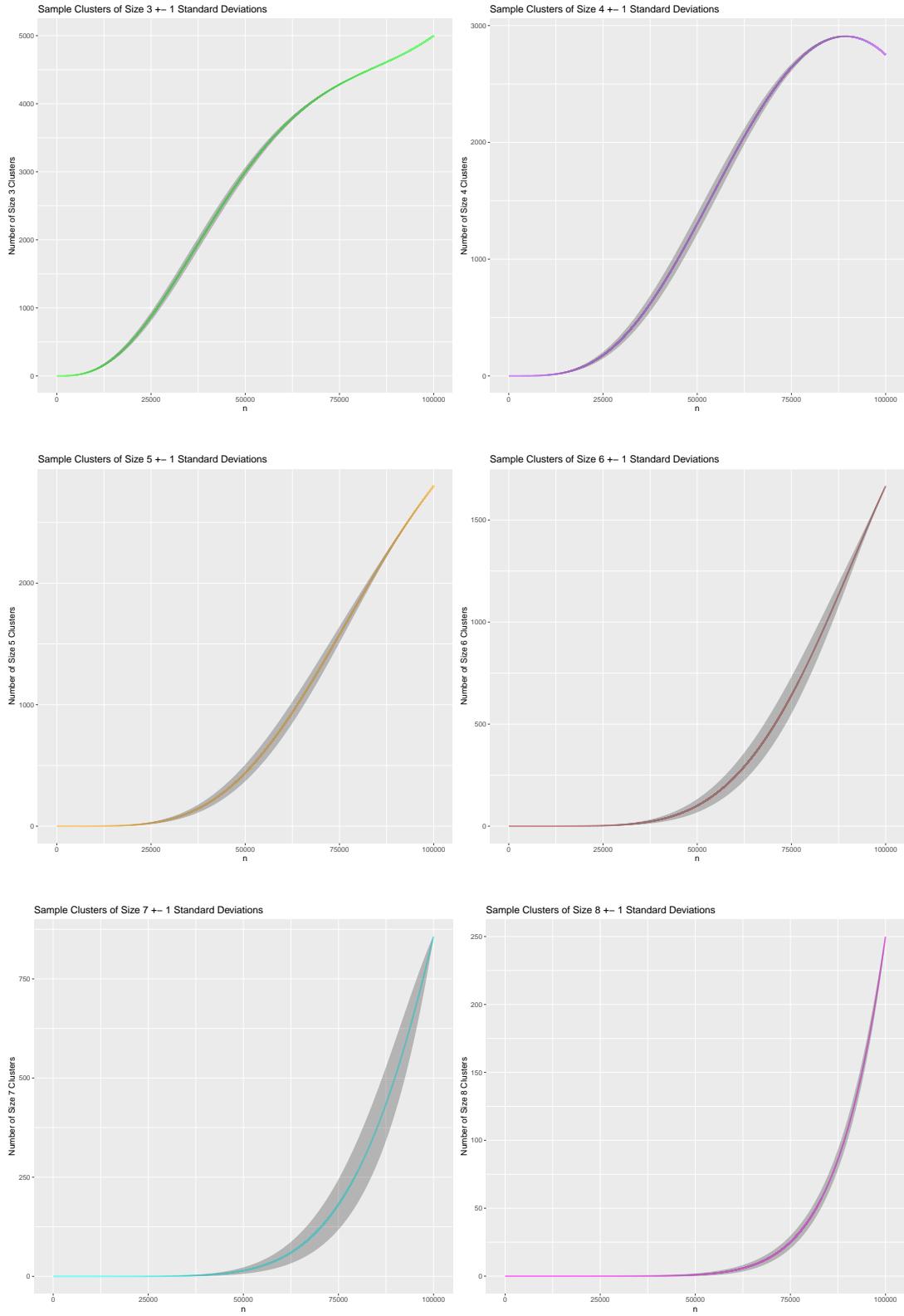


Figure 3: Expectations \pm 1 Standard Deviation for $c = 3, 4, 5, 6, 7, 8$.

These 8 plots show the number of sample clusters almost universally increasing as samples are taken. This consistent increase is due to the broadly decreasing distribution of cluster sizes (9). The exception is then seen for clusters of size 4, as there are more clusters of size 5 in the population. This creates the turning point towards the end of the process when the final molecule from clusters of size 5 are sampled and hence moving from size 4 to size 5 in the sample.

Additionally, we see that the standard deviation bounds are relatively small throughout with larger variances observed towards the middle and end of the process, as we might expect. We then observe the covariances between these clusters:

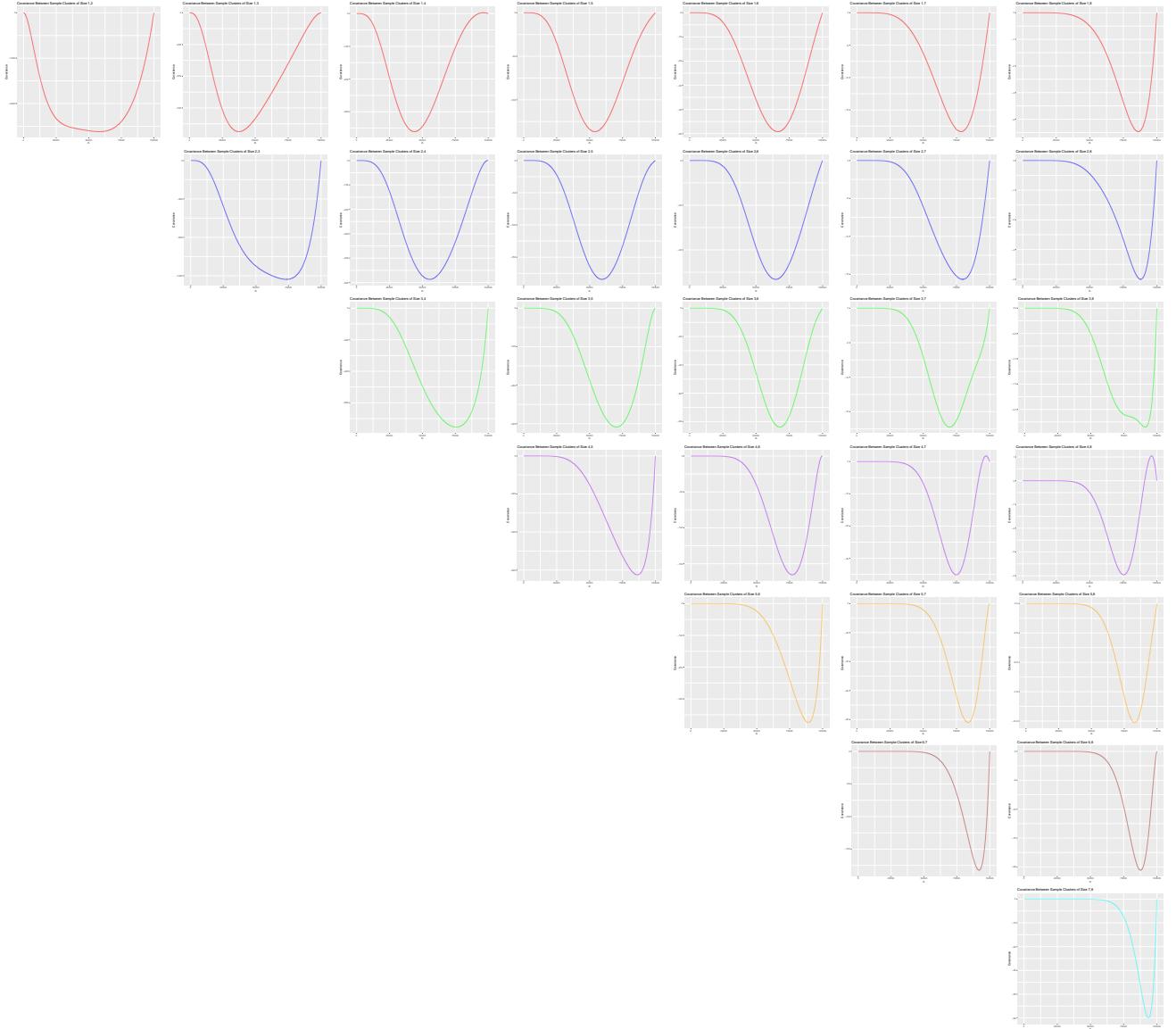
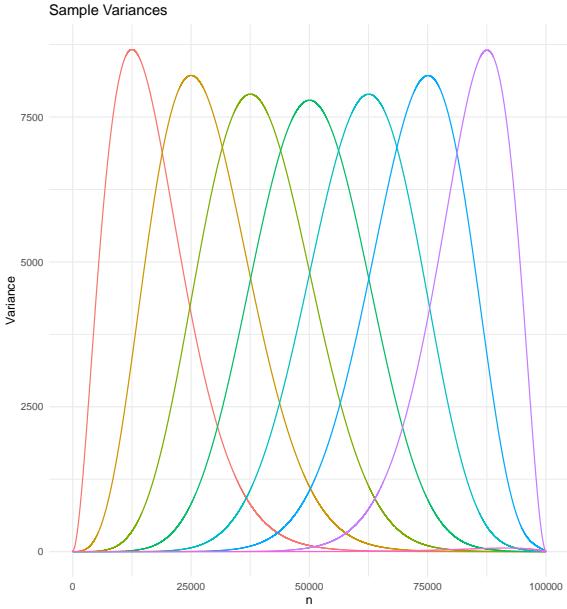


Figure 4: Covariance Plots.

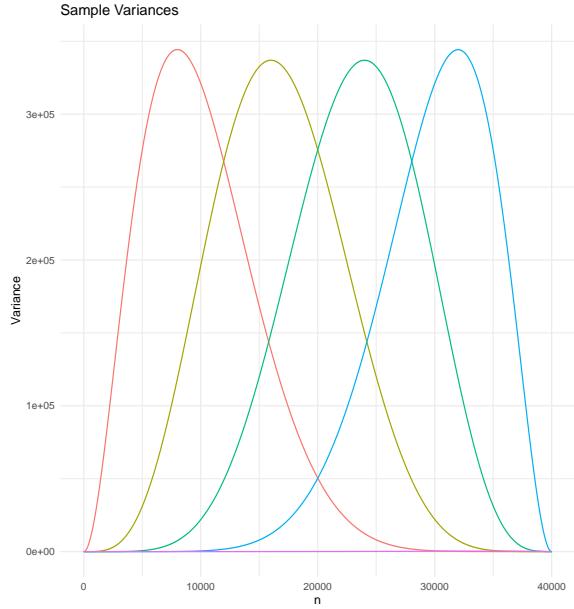
These covariances are broadly negative which we may expect. For example, if a sample cluster of size 5 is formed, it has sequentially increased from a size 1 or singleton cluster. So for $S_{1,5}(n)$ to increase, with some delay, $S_{1,1}(n), \dots, S_{1,4}(n)$ all have to decrease.

Exceptions to this behaviour include the relationship between clusters of size 4, and clusters of size 7 and 8. As the process nears completion, a small positive relationship can be seen. The reason for this is unclear, however, due to its small amplitude, it may reflect the small values of $\alpha_{1,7}$ and $\alpha_{1,8}$. As there are so few large clusters, it could be that the chance of even 4 molecules being sampled from them early in the process is slim, and thus these clusters are rapidly sampled from in the final stage of sampling. This could then correspond to a brief increase in the number of size 4 clusters as size 7 and 8 clusters start to form.

Finally we present the sample variances over the process. Regardless of the initial distribution, a remarkable symmetry exists between the variances of different sample cluster sizes.



(a) Sample Variances for the Described Simulation



(b) Sample Variances for $M=5$ with $\alpha_{1,c} = 0.2$

Whilst we may expect symmetry under the homogeneous initial distribution on the right where size-biased behaviour is nullified, seeing this across all scenarios is surprising. These graphs lead us to conjecture the following:

$$\text{Var}(S_{q,c}(n)) = \text{Var}(S_{q,M_q-c}(N-n)), \text{ for } c = 1, \dots, M_q - 1$$

3.7 Population Results

The second moment, variance, and covariance are determined in the same manner for the population as for the sample, however, they are of less interest in the context of the dilution process and seem to attract no consideration across the literature.

Theorem 3.1. *The population results are given as:*

$$\begin{aligned} \mathbb{E}[X_{q,c}(n)^2] &= \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{\binom{k}{k-c} \binom{N-k}{n-k+c}}{\binom{N}{n}} - \frac{\binom{k}{k-c} \binom{k}{k-c} \binom{N-2k}{n-2k+2c}}{\binom{N}{n}} \right] \\ &\quad + \sum_{k=c}^{M_q} \sum_{l=c}^{M_q} \left[\frac{\binom{k}{k-c} \binom{l}{l-c} \binom{N-k-l}{n-k-l+2c}}{\binom{N}{n}} \right] \\ &= \mathbb{E}[S_{q,c}(N-n)^2] \end{aligned}$$

so as an immediate consequence ,

$$\text{Var}(X_{q,c}(n)) = \text{Var}(S_{q,c}(N-n))$$

Likewise for $c_1 < c_2$ we find :

$$\text{Cov}(X_{q_1,c_1}(n), X_{q_2,c_2}(n)) = \text{Cov}(S_{q_1,c_1}(N-n), S_{q_2,c_2}(N-n))$$

Proof. See 8.2. □

The symmetry between population and sample can be loosely explained using the indicator random variables from which they are composed. Ideally we wish to explain a relationship between:

$$\mathbb{I}_{i_k,k,q}^{k-c}(n), \text{ and } \mathbb{I}_{i_k,k,q}^c(N-n)$$

It must be that all k molecules of the i_k 'th cluster are sampled in N steps. Additionally, by viewing the process ‘in reverse’ over all N steps, it would appear as a normal iteration of the full dilution process.

Therefore, if $k-c$ molecules were sampled in the forward direction in n steps, this directly corresponds to sampling c times from this cluster in $N-n$ steps in the reverse direction, or, directly to sampling c times in the remaining $N-n$ steps of the forward direction. This argument could be made rigorous using a coupling argument; however, we do not pursue this here.

3.8 Average Cluster Size

As well as knowing the expected number of each cluster size in the sample, it is of interest to know the expected average cluster size after a given number of samples. In the context of dilution in PCR, clusters of a specific size are optimal for inference, and thus we would like to know after how many samples this occurs. Letting:

$$\mu_q(n) = \frac{\sum_{c=1}^{M_q} c \cdot S_{q,c}(n)}{\sum_{c=1}^{M_q} S_{q,c}(n)} \quad (10)$$

be the average cluster size of type q molecules in the sample, by definition it is the total number of type q molecules divided by the total number of clusters. We note that in the case $Q = 1$, the numerator reduces to n .

Whilst both the numerator and denominator are easy to derive moments for, it is generally not true that:

$$\mathbb{E}\left[\frac{X}{Y}\right] = \frac{\mathbb{E}[X]}{\mathbb{E}[Y]}$$

for random variables, deriving information such as $\mathbb{E}[\mu_q(n)]$ is difficult, and therefore we use the ratio of expectations to describe the average cluster size. Defining the numerator of (10) as $N_q(n)$, the number of molecules of type q sampled in n samples, we have by linearity:

$$\begin{aligned} \mathbb{E}[N_q(n)] &= \mathbb{E}\left[\sum_{c=1}^M c \cdot S_{q,c}(n)\right] \\ &= \sum_{c=1}^M \sum_{k=c}^M c X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} \\ &= \sum_{k=1}^M X_{q,k}(0) \sum_{c=1}^k c \cdot \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} \\ &= \sum_{k=1}^M X_{q,k}(0) \frac{nk}{N} = \sum_{k=1}^M \alpha_{q,k} n \end{aligned} \quad (11)$$

where to obtain (11), we have recognised the expectation of a hypergeometric random variable, then used our earlier definition of $\alpha_{q,k}$. This simple closed-form expression alludes to the intuition that $N_q(n)$ should itself be hypergeometrically distributed by considering type q molecules as the type 1

objects being sampled without replacement. We then deduce that:

$$N_q(n) \sim \text{Hypergeometric}(N, \sum_{k=1}^{M_q} \alpha_{q,k} N, n)$$

and as an immediate result:

$$\text{Var}(N_q(n)) = \left(\sum_{k=1}^{M_q} \alpha_{q,k} n \right) \left(1 - \sum_{k=1}^{M_q} \alpha_{q,k} \right) \frac{N-n}{N-1}$$

We denote the denominator of (10) which describes the total number of clusters in the sample of type q by $C_q(n)$. Returning to the indicator definition of $S_{q,c}(n)$, $C_q(n)$ is defined as:

$$\begin{aligned} C_q(n) &= \sum_{c=1}^{M_q} \sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k, k, q}^c(n) \\ &= \sum_{k=1}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \sum_{c=1}^k \mathbb{1}_{i_k, k, q}^c(n) \\ &= \sum_{k=1}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} (1 - \mathbb{1}_{i_k, k, q}^0(n)) \end{aligned}$$

where $\mathbb{1}_{i_k, k, q}^0(n)$ takes on the value 1 if no samples have been taken from cluster i_k . Clearly we are then summing the number of original clusters from which at least one sample has been taken. From this, we compute:

$$\mathbb{E}[C_q(n)] = \sum_{k=1}^{M_q} X_{q,k}(0) \left[1 - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right] \quad (12)$$

which we will see later on is a more natural random variable associated with the sampling process of grouped elements.

3.8.1 Simulation Results

For the same simulation parameters described in 3.6, we compare the expectation of ratios given by:

$$\hat{\mu}_n = \frac{n}{\sum_{k=1}^{M_q} X_{q,k}(0) \left[1 - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right]} \quad (13)$$

with the true average cluster size for 3 sample paths of the process.

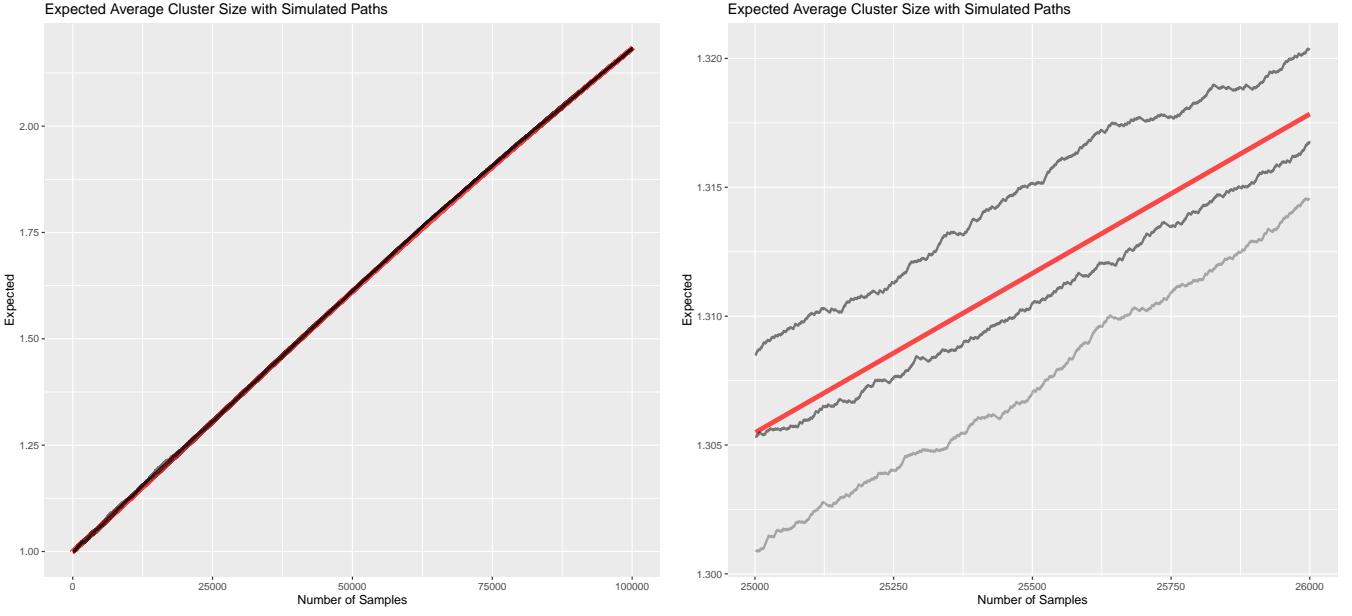


Figure 6: Average Cluster Size over 100,000 Samples. The Red Line Represents our Estimate

Over the entire process, we see that the sample paths, given as grey lines, match the ratio of expectations exceptionally well. Looking at a closer scale, we see that the true average cluster size falls within 0.005 of our estimate for these paths. In future, it would be desirable to estimate a corresponding variance, especially in the simpler case of $Q = 1$, however, we have been unable to achieve this.

4 Sampling with Replacement

Sampling without replacement is inherent to the dilution process and thus, considering sampling with replacement has no physical value. However, it is important to consider sampling with replacement from the population as described in 2.2 not only for potential use in other contexts, but also as the results of sampling with and without replacement typically converge for a large enough population.

When sampling from a clustered population with replacement, as before, we group molecules in the sample with those belonging to the same original cluster. Additionally, there now exists the possibility of sampling molecules from a cluster of size c more than c times. Therefore, using \tilde{S} to discuss sampling with replacement, we find:

$$\mathbb{E}[\tilde{S}_{q,c}(n)] = \sum_{k=1}^{M_q} X_{q,k}(0) \mathbb{P}(\text{Sampling } c \text{ times from a size } k \text{ cluster})$$

Whereas when sampling without replacement such probabilities were hypergeometric, when sampling with replacement they are binomial. Therefore;

$$\mathbb{E}[\tilde{S}_{q,c}(n)] = \sum_{k=1}^{M_q} X_{q,k}(0) \binom{n}{c} \left(\frac{k}{N}\right)^c \left(1 - \frac{k}{N}\right)^{n-c}$$

4.0.1 Derivation of Higher Moments

The variance and covariance of $\tilde{S}_{q,c}$ can be found using indicators in exactly the same manner as before. As mentioned above, the key difference is that the sum is now defined over all possible cluster

sizes:

$$\tilde{S}_{q,c}(n) = \sum_{k=1}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,q,k}^c(n)$$

The second difference is that when computing the second moment $\mathbb{E}[\tilde{S}_{q,c}(n)^2]$, and hence the expectation of indicator products such as:

$$\mathbb{1}_{i_k,q,k}^c(n) \mathbb{1}_{i_l,q,l}^c(n) = \begin{cases} 1, & \text{if we sample } c \text{ times from a size } k \text{ cluster and from a size } l \\ & \text{cluster in } n \text{ samples with replacement} \\ 0, & \text{otherwise} \end{cases}$$

Whereas we found $\mathbb{P}(\mathbb{1}_{i_k,q,k}^c(n) \mathbb{1}_{i_l,q,l}^c(n) = 1)$ to be multivariate hypergeometric when sampling without replacement, now, we naturally find this to be multinomial. These observations lead us to the following.

Theorem 4.1. *For a population of molecules of size N , where each molecule is of type q and belongs to a cluster or group of size c for $c = 1, \dots, M_q$, the number of size c clusters of type q molecules after n samples with replacement satisfy:*

$$\begin{aligned} \mathbb{E}[\tilde{S}_{q,c}(n)^2] &= \sum_{k=1}^{M_q} X_{q,k}(0) \left[\binom{n}{c} \left(\frac{k}{N}\right)^c \left(1 - \frac{k}{N}\right)^{n-c} - \binom{n}{c, c, n-c} \left(\frac{k}{N}\right)^{2c} \left(1 - \frac{2k}{N}\right)^{n-2c} \right] \\ &\quad + \sum_{k=1}^{M_q} \sum_{l=1}^{M_q} X_{q,k}(0) X_{q,l}(0) \left[\binom{n}{c, c, n-2c} \left(\frac{k}{N}\right)^c \left(\frac{l}{N}\right)^c \left(1 - \frac{k+l}{N}\right)^{n-c} \right] \end{aligned}$$

and as a consequence:

$$\begin{aligned} \text{Var}(\tilde{S}_{q,c}(n)) &= \sum_{k=1}^{M_q} \sum_{l=1}^{M_q} X_{q,k}(0) X_{q,l}(0) \binom{n}{c, c, n-2c} \left(\frac{k}{N}\right)^c \left(\frac{l}{N}\right)^c \left(1 - \frac{k+l}{N}\right)^{n-c} \\ &\quad - \sum_{k=1}^{M_q} \sum_{l=1}^{M_q} X_{q,k}(0) X_{q,l}(0) \binom{n}{c}^2 \left(\frac{k}{N}\right)^c \left(1 - \frac{k}{N}\right)^{n-c} \left(\frac{l}{N}\right)^c \left(1 - \frac{l}{N}\right)^{n-c} \\ &\quad + \sum_{k=1}^{M_q} X_{q,k}(0) \left[\binom{n}{c} \left(\frac{k}{N}\right)^c \left(1 - \frac{k}{N}\right)^{n-c} - \binom{n}{c, c, n-2c} \left(\frac{k}{N}\right)^{2c} \left(1 - \frac{2k}{N}\right)^{n-2c} \right] \end{aligned}$$

Furthermore, for distinct sizes $c_1 < c_2$ and any two molecule types q_1, q_2 :

$$\begin{aligned} \text{Cov}(\tilde{S}_{q_1,c_1}(n), \tilde{S}_{q_2,c_2}(n)) &= \sum_{k=1}^{M_{q_1}} \sum_{l=1}^{M_{q_2}} X_{q_1,k}(0) X_{q_2,l}(0) \binom{n}{c_1, c_2, n - c_1 - c_2} \left(\frac{k}{N}\right)^{c_1} \left(\frac{l}{N}\right)^{c_2} \left(1 - \frac{k+l}{N}\right)^{n-c_1-c_2} \\ &\quad - \sum_{k=1}^{M_{q_1}} \sum_{l=1}^{M_{q_2}} X_{q_1,k}(0) X_{q_2,l}(0) \binom{n}{c_1} \binom{n}{c_2} \left(\frac{k}{N}\right)^{c_1} \left(1 - \frac{k}{N}\right)^{n-c_1} \left(\frac{l}{N}\right)^{c_2} \left(1 - \frac{l}{N}\right)^{n-c_2} \\ &\quad - \delta_{q_1, q_2} \sum_{k=1}^{M_q} X_{q_1,k}(0) \binom{n}{c_1, c_2, n - c_1 - c_2} \left(\frac{k}{N}\right)^{c_1+c_2} \left(1 - \frac{2k}{N}\right)^{n-c_1-c_2} \end{aligned}$$

These moments mirror exactly the structure of results obtained when sampling without replacement.

5 Species Richness Estimation

At the core of the dilution process was estimating the number of clusters of a specific size when sampling without replacement from a population composed of varied groups. If we cease to group clusters by their sizes in both the initial population and sample, the population of size N is divided into S groups, the total number of clusters.

$$S = \sum_{q=1}^Q \sum_{c=1}^{M_q} X_{q,c}(0)$$

Typically, the above situation describes a population of S species of animals, in which case, previous sections correspond to estimating the number of species represented by c specimens in a sample of size n .

The species are of size N_i for $i = 1, \dots, S$ and we denote the number of distinct species observed in a sample without replacement of size n by \hat{S}_n . Letting X_i be the number of individuals sampled from species i , then clearly,

$$\hat{S}_n = \sum_{i=1}^S \mathbb{1}_{\{X_i > 0\}}(n)$$

where:

$$\mathbb{1}_{\{X_i > 0\}}(n) = \begin{cases} 1, & \text{if the } i\text{'th species is observed at all in } n \text{ samples} \\ 0, & \text{if the } i\text{'th species is unobserved in } n \text{ samples} \end{cases}$$

in which case,

$$\begin{aligned} \mathbb{E}[\hat{S}_n] &= \sum_{i=1}^S \mathbb{P}(\mathbb{1}_{\{X_i > 0\}}(n) = 1) \\ &= \sum_{i=1}^S [1 - \mathbb{P}(X_i = 0)] \\ &= \sum_{i=1}^S \left[1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right] \end{aligned}$$

Which is stated but not derived in [7]. We also note that by letting $S(k)$ be the number of distinct species of size k in the population, so that $\sum_{k=1}^M S(k) = S$, where M is the largest species size, we can alternatively define the expected value as;

$$\mathbb{E}[\hat{S}_n] = \sum_{k=1}^M S(k) \left[1 - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right]$$

which we recognise as the total cluster number (12) derived earlier when looking for the average cluster size. Note that this corresponds to;

$$S(k) = \sum_{q=1}^Q X_{q,k}(0)$$

We infer from this that \hat{S}_n can be decomposed into $\hat{S}_n(c)$ of the form described earlier.

The variance of the overall species estimator is given in [8] and derived in 8.3 as;

$$\text{Var}(\hat{S}_n) = \sum_{i=1}^S \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right) + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \left(\frac{\binom{N-N_i-N_j}{n}}{\binom{N}{n}} - \frac{\binom{N-N_i}{n} \binom{N-N_j}{n}}{\binom{N}{n}^2} \right)$$

or alternatively,

$$\text{Var}(\hat{S}_n) = \sum_{k=1}^M S(k) \left(\frac{\binom{N-k}{n}}{\binom{N}{n}} - \frac{\binom{N-2k}{n}}{\binom{N}{n}} \right) + \sum_{k=1}^M \sum_{l=1}^M S(k) S(l) \left(\frac{\binom{N-k-l}{n}}{\binom{N}{n}} - \frac{\binom{N-k}{n} \binom{N-l}{n}}{\binom{N}{n}^2} \right)$$

Likewise, in the case that sampling is conducted with replacement, we denote the overall number of species observed as \tilde{S}_n and find that:

$$\begin{aligned} \mathbb{E}[\tilde{S}_n] &= \sum_{i=1}^S (1 - (1 - \frac{N_i}{N})^n) \\ \text{Var}(\tilde{S}_n) &= \sum_{i=1}^S \left[(1 - \frac{N_i}{N})^n - (1 - \frac{2N_i}{N})^n \right] + \sum_{i=1}^S \sum_{j=1}^S \left[(1 - \frac{N_i+N_j}{N})^n - (1 - \frac{N_i}{N})^n (1 - \frac{N_j}{N})^n \right] \\ &= \sum_{k=1}^M S(k) \left[(1 - \frac{k}{N})^n - (1 - \frac{2k}{N})^n \right] + \sum_{k=1}^M \sum_{l=1}^M S(k) S(l) \left[(1 - \frac{k+l}{N})^n - (1 - \frac{k}{N})^n (1 - \frac{l}{N})^n \right] \end{aligned}$$

which as expected, displays the same structure as the without replacement case. We also note that now the clusters remaining in the population correspond to the number of unobserved species and are therefore given by $S - \hat{S}$ or $S - \tilde{S}$ in each case.

6 Asymptotic Results

We wish to consider the asymptotic behaviour of our estimators $S_{q,c}(n)$ and \hat{S}_n as $N \rightarrow \infty$. Determining these results will help us uncover when these estimators are asymptotically unbiased for the number of species.

Instead of using a fixed sample size, we will also consider $n \rightarrow \infty$, however, we will need to account for the fact this can occur in various ways relative to N . For example, we could have $n = o(\sqrt{N})$ or $\frac{n}{N} \rightarrow p \in (0, 1)$ as $n, N \rightarrow \infty$.

Additionally, there is a diverse range of scenarios for the initial population that we should consider. On either end of this range we could have:

- A population composed only of singleton species as $N \rightarrow \infty$, i.e. $S(1) = N$,
- A population dominated by a small number M of large species so that $N_i \rightarrow \infty$ for $i = 1, \dots, M$ as $N \rightarrow \infty$

As a third example, we could have a combination of these cases wherein large species dominate half of the overall population, and singletons make up the remaining size.

We begin by addressing an asymptotic result from the literature. Translating the notation of [6] into our own, the following is given.

Theorem 6.1. If as $n, N \rightarrow \infty$, $\frac{n}{N} \rightarrow p \in (0, 1)$, then:

$$\mathbb{E}[S_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \binom{k}{c} p^c (1-p)^{k-c} (1 + o(1))$$

No proof of this result is provided, however, we give our own.

Proof. Stirling's approximation returns:

$$n! \sim \sqrt{2\pi n} n^n e^{-n}, \quad n \rightarrow \infty$$

i.e. the error is $o(1)$. Therefore,

$$\begin{aligned} \frac{\binom{N-k}{n-c}}{\binom{N}{n}} &= \frac{(N-k)!n!(N-n)!}{(N-k-n+c)!(n-c)!N!} \\ &\sim \sqrt{\frac{(N-k)n(N-n)}{(N-k-n+c)(n-c)N}} \cdot \frac{(N-k)^{N-k} n^n (N-n)^{N-n}}{(N-k-n+c)^{N-k-n+c} (n-c)^{n-c} N^N} \end{aligned}$$

Now, we separate the terms into finite powers of k, c and use that;

$$\sqrt{\frac{(N-k)n(N-n)}{(N-k-n+c)(n-c)N}} \rightarrow \sqrt{\frac{N \cdot n \cdot (N-n)}{(N-n) \cdot n \cdot N}} = 1, \quad n, N \rightarrow \infty$$

Hence, we find:

$$\begin{aligned} \frac{\binom{N-k}{n-c}}{\binom{N}{n}} &\sim \left(\frac{N-n}{N-k-n+c} \right)^{N-n} \left(\frac{n}{n-c} \right)^n \left(\frac{N-k}{N} \right)^N \cdot \frac{(N-k)^{-k}}{(N-k-n+c)^{-k+c} (n-c)^{-c}} \\ &= \left(1 - \frac{k-c}{N-n} \right)^{-(N-n)} \left(1 - \frac{c}{n} \right)^{-n} \left(1 - \frac{k}{N} \right)^N \cdot \left(\frac{N-k-n+c}{N-k} \right)^{k-c} \cdot \left(\frac{n-c}{N-k} \right)^c \end{aligned}$$

Finally, using the definition of the exponential function:

$$\lim_{n \rightarrow \infty} \left(1 - \frac{x}{n} \right)^n = e^{-x}$$

We find that as $n, N \rightarrow \infty$:

$$\left(1 - \frac{k-c}{N-n} \right)^{-(N-n)} \left(1 - \frac{c}{n} \right)^{-n} \left(1 - \frac{k}{N} \right)^N \rightarrow e^{k-c} \cdot e^c \cdot e^{-k} = 1$$

So that we ultimately have:

$$\frac{\binom{N-k}{n-c}}{\binom{N}{n}} \sim \left(\frac{n}{N} \right)^c \left(1 - \frac{n}{N} \right)^{k-c}$$

As this holds for all $k \geq c$, we hence find that as $n, N \rightarrow \infty$:

$$\mathbb{E}[S_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-c}}{\binom{N}{n}} = \sum_{k=c}^{M_q} X_{q,k}(0) \binom{k}{c} p^c (1-p)^{k-c} (1 + o(1))$$

□

6.1 Extra Notation

Moving forward, we introduce the random variable ξ describing the size of a randomly chosen cluster i.e.

$$\mathbb{P}(\xi = k) = \frac{S(k)}{S}, \text{ for } k = 1, \dots, M$$

where for shorthand we also use $\psi(k) = \mathbb{P}(\xi = k)$. Note that:

$$\mathbb{E}[\xi] = \sum_{k=1}^M \frac{kS(k)}{S} = \frac{1}{S} \sum_{k=1}^M kS(k) = \frac{N}{S}$$

is the average cluster size, and that

$$\text{Var}(\xi) = \sum_{k=1}^M \frac{k^2 S(k)}{S} - \frac{N^2}{S^2}$$

We index all relevant variables by n to indicate their asymptotic dependence and introduce 3 new variables.

$$\hat{\nu}_n = \frac{\hat{S}_n}{S_n}, \quad \tilde{\nu}_n = \frac{\tilde{S}_n}{S_n}, \quad \lambda_n = \frac{n}{N}, \quad \xi_n \text{ and } M_n$$

$\hat{\nu}_n$ describes the proportion of species observed without replacement, $\tilde{\nu}_n$, for sampling with replacement, and lastly, λ_n is the proportion of the population sampled. We also trivially require that:

$$\limsup_{n \rightarrow \infty} \lambda_n < 1, \quad \text{and} \quad \limsup_{n \rightarrow \infty} \frac{M_n}{N} < 1$$

With this new notation, we can rephrase Shlosser's result 6.1 as:

$$\frac{\mathbb{E}[S_{q,c}(n)]}{S_n} = \sum_{k=c}^{M_q} \psi(k) p^c (1-p)^{k-c} (1 + o(1))$$

or,

$$\mathbb{E}[S_{q,c}(n)] = \mathbb{E}[\tilde{S}_{q,c}(n)] + o(1)$$

if $\lambda_n \rightarrow p$ as $n, N \rightarrow \infty$. i.e. the with and without replacement expectations coincide asymptotically.

6.2 Limit Theorems

Theorem 6.2. As $n \rightarrow \infty$, assume that:

$$\lambda_n^2 \mathbb{E}[\xi_n e^{-\lambda_n \xi_n}] = o(1)$$

if this holds, then:

$$\mathbb{E}[\hat{\nu}_n] = \mathbb{E}[\tilde{\nu}_n] + o(1)$$

i.e. under suitable conditions, the expected without replacement proportion approaches the same result when sampling with replacement.

Proof. Firstly, note that in this approximation we are looking for a bound on the difference between these terms which are given by:

$$\mathbb{E}[\hat{\nu}_n] = 1 - \sum_{k=1}^M \psi(k) \frac{\binom{N-k}{n}}{\binom{N}{n}}$$

and:

$$\mathbb{E}[\tilde{\nu}_n] = 1 - \sum_{k=1}^M \psi(k) \left(1 - \frac{k}{N}\right)^n$$

so that we wish to bound:

$$|\mathbb{E}[\tilde{\nu}_n] - \mathbb{E}[\hat{\nu}_n]| = \sum_{k=1}^M \psi(k) \left| \left(1 - \frac{k}{N}\right)^n - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right|$$

or, letting $Q = N - k$, and $q = \frac{Q}{N}$ we must bound:

$$\left| q^n - \frac{\binom{Q}{n}}{\binom{N}{n}} \right|$$

by definition of the binomial coefficient we have:

$$\frac{\binom{Q}{n}}{\binom{N}{n}} = \frac{Q!(N-n)!}{(Q-n)!N!} = \frac{Q \cdots (Q-n+1)}{N \cdots (N-n+1)} = \prod_{k=0}^{n-1} \frac{Q-k}{N-k}$$

Using a telescopic sum, it follows that:

$$q^n - \prod_{k=0}^{n-1} \frac{Q-k}{N-k} = \sum_{j=1}^n \left(q^j \prod_{k=j}^{n-1} \frac{Q-k}{N-k} - q^{j-1} \prod_{k=j-1}^{n-1} \frac{Q-k}{N-k} \right)$$

where we take:

$$\prod_{k=n}^{n-1} \frac{Q-k}{N-k} = 1$$

Taking out a common factor within the telescoping sum we have:

$$q^n - \prod_{k=0}^{n-1} \frac{Q-k}{N-k} = \sum_{j=1}^n \left(q^{j-1} \prod_{k=j}^{n-1} \frac{Q-k}{N-k} \left[\frac{Q}{N} - \frac{Q-j+1}{N-j+1} \right] \right)$$

As $q < 1$, we have that:

$$\frac{Q-k}{N-k} < \frac{Q}{N}$$

for each value of k in the sum. Additionally,

$$\frac{Q}{N} - \frac{Q-j+1}{N-j+1} = \frac{(N-Q)(j-1)}{N(N-j+1)} \leq \frac{(N-Q)(j-1)}{N(N-n)}$$

where this follows as $j \leq n-1$ and hence the denominator is minimised for $j = n-1$. Applying these we find that:

$$\begin{aligned} q^n - \prod_{k=0}^{n-1} \frac{Q-k}{N-k} &\leq \sum_{j=1}^n \left(\frac{(N-Q)(j-1)}{N(N-n)} q^{j-1+n-j} \right) \\ &= \frac{(N-Q)}{N(N-n)} \cdot q^{n-1} \sum_{j=1}^n (j-1) \\ &= (1-q) \cdot q^{n-1} \cdot \frac{n(n-1)}{2(N-n)} \end{aligned}$$

Hence, returning to the original terms, with $Q = N - k$ we have:

$$\begin{aligned}
|\mathbb{E}[\tilde{\nu}_n] - \mathbb{E}[\hat{\nu}_n]| &\leq \sum_{k=1}^M \psi(k) \frac{k}{N} (1 - \frac{k}{N})^{n-1} \frac{n(n-1)}{2(N-n)} \\
&\leq \sum_{k=1}^M k \psi(k) (1 - \frac{k}{N})^n \cdot \frac{N}{N(N-k)} \cdot \frac{n(n-1)}{(N-n)} \\
&\leq \sum_{k=1}^M k \psi(k) e^{-\frac{nk}{N}} \frac{n(n-1)}{2(N-n)(N-k)} \\
&\leq \sum_{k=1}^M k \psi(k) e^{-\frac{nk}{N}} \frac{n^2}{N^2} \\
&= \lambda_n^2 \mathbb{E}[\xi_n e^{-\xi_n \lambda_n}]
\end{aligned} \tag{14}$$

Which by the stated condition is $o(1)$ as $n, N \rightarrow \infty$.

Note: For (14) we have used the following inequality which is true for $m \geq 1$ and $|x| \leq m$

$$(1 - \frac{x}{m})^m \leq e^{-x}$$

Hence, we have that:

$$(1 - \frac{kn}{Nn})^n \leq e^{-\frac{kn}{N}}$$

□

Hence, under certain conditions the with replacement and without replacement estimates coincide.

Theorem 6.3. As $n \rightarrow \infty$, if $\mathbb{E}[\xi_n^2] = o(N)$, then:

$$\mathbb{E}[\hat{\nu}_n] = 1 - \mathbb{E}[(1 - \lambda_n)^{\xi_n}] + o(1)$$

Proof. In light of these two terms, we are looking to bound:

$$\begin{aligned}
|\mathbb{E}[\hat{\nu}_n] - (1 - \mathbb{E}[(1 - \lambda_n)^{\xi_n}])| &= \left| 1 - \sum_{k=1}^M \psi(k) \frac{\binom{N-k}{n}}{\binom{N}{n}} - \left(1 - \sum_{k=1}^M \psi(k) (1 - \frac{n}{N})^k \right) \right| \\
&= \sum_{k=1}^M \psi(k) \left| (1 - \frac{n}{N})^k - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right|
\end{aligned}$$

which is similar to the previous term, albeit with the roles of n and k reversed in the approximation. Given the stated assumption, we seek to bound this by:

$$C \sum_{k=1}^M \frac{k^2 \psi(k)}{N}$$

for some constant C .

We may re-express the hypergeometric probability as;

$$\frac{\binom{N-k}{n}}{\binom{N}{n}} = \prod_{i=0}^{n-1} \frac{N-k-i}{N-i} = \prod_{i=0}^{n-1} \left(1 - \frac{k}{N-i}\right) = \exp \left(\sum_{i=0}^{n-1} \ln \left(1 - \frac{k}{N-i}\right) \right)$$

We note that as $\frac{k}{N-i} > 0$, $\sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) < 0$. Likewise we can express the approximation using the exponential function:

$$(1 - \frac{n}{N})^k = \exp(k \ln(1 - \frac{n}{N}))$$

where $\ln(1 - \frac{n}{N}) < 0$. It is known for non-positive x and y that:

$$|\exp(x) - \exp(y)| < |x - y|$$

so it suffices for us to find a bound on:

$$\left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) - k \ln(1 - \frac{n}{N}) \right|$$

By the inequality:

$$|\ln(1 - x) + x| \leq \frac{x^2}{1 - x}$$

we have that:

$$|\ln(1 - \frac{k}{N-i}) + \frac{k}{N-i}| \leq \frac{k^2}{(N-i)(N-i-k)}$$

so that

$$\begin{aligned} & \left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) + k \sum_{i=0}^{n-1} \frac{1}{N-i} \right| \\ &= \left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) + k H_{N-n+1}^N \right| \\ &\leq \sum_{i=0}^{n-1} \frac{k^2}{(N-i)(N-i-k)} \end{aligned}$$

where;

$$H_m^n = \sum_{i=m}^n \frac{1}{i}$$

is the partial harmonic sum. Therefore, using that the denominator is decreasing in i , we find:

$$\left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) + k H_{N-n+1}^N \right| \leq \frac{k^2 n}{(N-n)(N-n-k)} \leq \frac{A k^2}{N}$$

for some constant A . Then, using the triangle inequality we therefore have:

$$\begin{aligned} & \left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) - k \ln(1 - \frac{n}{N}) \right| \\ &= \left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) + k H_{N-n+1}^N - k H_{N-n+1}^N - k \ln(1 - \frac{n}{N}) \right| \\ &\leq \left| \sum_{i=0}^{n-1} \ln(1 - \frac{k}{N-i}) + k H_{N-n+1}^N \right| + k |H_{N-n+1}^N + \ln(1 - \frac{n}{N})| \end{aligned}$$

To bound this final term, we note that for the Harmonic sum, we have the bounds (see 8.4):

$$\ln\left(\frac{n+1}{m}\right) \leq H_m^n \leq \ln\left(\frac{n}{m-1}\right)$$

so that in our case,

$$\ln\left(\frac{N+1}{N-n+1}\right) \leq H_{N-n+1}^N \leq \ln\left(\frac{N}{N-n}\right) = -\ln\left(1-\frac{n}{N}\right)$$

and hence:

$$\begin{aligned} |H_{N-n+1}^N + \ln(1 - \frac{n}{N})| &\leq \left| \ln\left(\frac{N+1}{N-n+1}\right) + \ln\left(1 - \frac{n}{N}\right) \right| \\ &\leq \left| \ln\left(\frac{(N+1)(N-n)}{N(N-n+1)}\right) \right| \\ &= \left| \ln\left(\frac{N(N-n+1)-n}{N(N-n+1)}\right) \right| \\ &= \left| \ln\left(1 - \frac{n}{N(N-n+1)}\right) \right| \end{aligned}$$

Using the approximation $\ln(1+x) \sim x$ for small x , we therefore find;

$$\begin{aligned} |H_{N-n+1}^N + \ln(1 - \frac{n}{N})| &\leq \left| -\frac{n}{N(N-n+1)} \right| \\ &= \frac{n}{N(N-n+1)} \\ &\leq \frac{B}{N} \end{aligned}$$

for some constant B . Combining the two derived bounds, we have that:

$$\begin{aligned} \sum_{k=1}^M \psi(k) \left| \left(1 - \frac{n}{N}\right)^k - \frac{\binom{N-k}{n}}{\binom{N}{n}} \right| &\leq \sum_{k=1}^M \psi(k) \left(\frac{Ak^2}{N} + \frac{Bk}{N} \right) \\ &\leq C \sum_{k=1}^M \frac{k^2 \psi(k)}{N} \end{aligned}$$

and hence by our condition, the asymptotic result holds. \square

Theorem 6.4. If $\mathbb{E}[\xi_n^2] = o(N)$ and $\lambda_n^2 \mathbb{E}[\xi e^{-\lambda_n \xi_n}] = o(1)$, then;

$$\mathbb{E}[\hat{\nu}_n] = 1 - \mathbb{E}[e^{-\lambda_n \xi_n}] + o(1)$$

Proof. We have already shown that under the first condition,

$$\mathbb{E}[\hat{\nu}_n] = 1 - \mathbb{E}[(1 - \lambda_n)^{\xi_n}] + o(1)$$

so need only show that under the additional second condition,

$$\mathbb{E}[(1 - \lambda_n)^{\xi_n}] \sim \mathbb{E}[e^{-\lambda_n \xi_n}]$$

The earlier inequality may be extended to:

$$(1 - \frac{x}{m})^m \leq e^{-x} \leq (1 - \frac{x}{m})^m + \frac{x^2}{m} e^{-x}$$

and hence, using $x = \frac{nk}{N}$, $m = k$ we find:

$$\left| \left(1 - \frac{n}{N}\right)^k - e^{-\frac{nk}{N}} \right| \leq \frac{n^2 k}{N^2} e^{-\frac{nk}{N}}$$

So finally, we have:

$$\begin{aligned}
|\mathbb{E}[(1 - \lambda_n)^{\xi_n}] - \mathbb{E}[e^{-\lambda_n \xi_n}]| &= \sum_{k=1}^M \psi(k) \left| (1 - \frac{n}{N})^k - e^{-\frac{nk}{N}} \right| \\
&\leq \lambda_n^2 \sum_{k=1}^M k \psi(k) e^{-k \lambda_n} \\
&= \lambda_n^2 \mathbb{E}[\xi e^{-\lambda_n \xi_n}]
\end{aligned}$$

so the result holds. \square

Therefore, we have valid asymptotics for the expected proportion of species observed without replacement given two separate conditions on the population, as well in the case of their intersection. In each case, these also suggest situations in which our estimate is asymptotically unbiased if the approximation is shown to approach 1.

The condition $\mathbb{E}[\xi_n^2] = o(N)$ seems to suggest suitability for lots of small species as in the extreme case of $\psi(1) = 1$, we have:

$$\mathbb{E}[\xi_n^2] = \sum_{k=1}^M k^2 \psi(k) = 1$$

whilst alternatively, if we suppose there are M large species of size $\frac{N}{M}$, we find;

$$\mathbb{E}[\xi_n^2] = \frac{N^2}{C^2} \neq o(N)$$

However, the condition $\lambda_n^2 \mathbb{E}[\xi e^{-\lambda_n \xi_n}]$ seems reasonable when larger dominant species are present as in the same scenario presented above:

$$\lambda_n^2 \mathbb{E}[\xi e^{-\lambda_n \xi_n}] = \frac{n^2}{N^2} \cdot \frac{N}{M} e^{-\frac{n}{N} \cdot \frac{N}{M}} = \frac{n}{N} \cdot \frac{n}{M} e^{-\frac{n}{M}} = o(1)$$

Whilst these represent extreme cases, they illustrate some physical intuition behind the conditions on our estimators.

7 Conclusion

In conclusion, the dilution process seems to exhibit a variety of rich and unexpected behaviour which warrants further investigation. Within one example, we have highlighted unexpected symmetries as well as a set of interesting covariance plots for which far more variations are sure to exist depending on the makeup of the population.

Additionally, as (13) seems to be a reasonable approximation to the average cluster size, it may be possible to use this result for improved computation time and easier inference within PCR, though, it would be desirable to have some idea of the error within this estimate. The asymptotic results shown may in theory be applied to a wide variety of possible scenarios, so it would be interesting to assess their performance on real datasets or further simulations.

In future research, it would also be desirable to investigate the general symmetries that exists between the population and sample when objects are drawn without replacement from a population composed of classes. Furthermore, finding an improved closed form estimate for the average cluster size would be of key benefit to current applications. Overall, our work highlights that despite the wide body of literature, there are still many interesting problems to be discovered within the broad topic of species richness estimation.

References

- [1] R. A. Fisher, A. S. Corbet, and C. B. Williams, “The relation between the number of species and the number of individuals in a random sample of an animal population,” *Journal of Animal Ecology*, vol. 12, no. 1, pp. 42–58, 1943. [Online]. Available: <http://www.jstor.org/stable/1411>
- [2] M. A. McGeoch, M. Schroeder, B. Ekbom, and S. Larsson, “Saproxylic beetle diversity in a managed boreal forest: importance of stand characteristics and forestry conservation measures,” *Diversity and Distributions*, vol. 13, no. 4, pp. 418–429, 2007. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1472-4642.2007.00350.x>
- [3] B. Efron and R. Thisted, “Estimating the number of unseen species: How many words did Shakespeare know?” *Biometrika*, vol. 63, no. 3, pp. 435–447, 1976. [Online]. Available: <http://www.jstor.org/stable/2335721>
- [4] J. Bunge and M. Fitzpatrick, “Estimating the number of species: A review,” *Journal of the American Statistical Association*, vol. 88, no. 421, pp. 364–373, 1993. [Online]. Available: <http://www.jstor.org/stable/2290733>
- [5] A. Chao and C.-H. Chiu, *Species Richness: Estimation and Comparison*. John Wiley Sons, Ltd, 2016, pp. 1–26. [Online]. Available: <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118445112.stat03432.pub2>
- [6] A. Shlosser, “On estimation of the size of a dictionary of a long text on the basis of a sample,” *Engineering Cybernetics*, vol. 19, pp. 97–102, 1981.
- [7] S. H. Hurlbert, “The nonconcept of species diversity: A critique and alternative parameters,” *Ecology*, vol. 52, no. 4, pp. 577–586, 1971. [Online]. Available: <http://www.jstor.org/stable/1934145>
- [8] K. L. Heck, G. van Belle, and D. Simberloff, “Explicit calculation of the rarefaction diversity measurement and the determination of sufficient sample size,” *Ecology*, vol. 56, no. 6, pp. 1459–1461, 1975. [Online]. Available: <http://www.jstor.org/stable/1934716>

8 Appendix

8.1 Recurrence Relations

First we should deal with the simplest case, $c = M_q$. We have:

$$X_{q,M_q}(n+1) = X_{q,M_q}(n) + \begin{cases} 0, & \text{if we sample from any other cluster size} \\ -1, & \text{if we sample from a cluster of size } M_q \end{cases}$$

or equivalently

$$X_{q,M_q}(n+1) = X_{q,M_q}(n) + \begin{cases} 0, & \text{with probability; } \frac{N-n-M_q X_{q,M_q}(n)}{N-n} \\ -1, & \text{with probability; } \frac{M_q X_{q,M_q}(n)}{N-n} \end{cases}$$

Taking first a conditional expectation with respect to $X_{q,M_q}(n)$, then taking the expectation again, we find:

$$\mathbb{E}(X_{q,M_q}(n+1)) = \frac{(N-n-M_q)\mathbb{E}(X_{q,M_q}(n))}{N-n} \tag{15}$$

in conjunction with initial condition $\mathbb{E}[X_{q,M_q}(0)] = X_{q,M_q}(0)$.

Lemma 8.1. *The recurrence relation (15) with corresponding initial condition is solved by;*

$$\mathbb{E}(X_{q,M_q}(n)) = X_{q,M_q}(0) \frac{\binom{N-M_q}{n}}{\binom{N}{n}} \quad (16)$$

Proof. Firstly, for $n = 0$ we have:

$$\mathbb{E}[X_{q,M_q}(0)] = X_{q,M_q}(0) \frac{\binom{N-M_q}{0}}{\binom{N}{0}} = X_{q,M_q}(0)$$

so this satisfies our initial condition. Then, plugging this into the right hand side of the recurrence relation we find:

$$\begin{aligned} &= X_{q,M_q}(0) \frac{N-n-M_q}{N-n} \cdot \frac{\binom{N-M_q}{n}}{\binom{N}{n}} \\ &= X_{q,M_q}(0) \frac{N-n-M_q}{N-n} \cdot \frac{(N-M_q)!(N-n)!n!}{(N-n-M_q)!n!N!} \cdot \frac{n+1}{n+1} \\ &= X_{q,M_q}(0) \frac{\binom{N-M_q}{[n+1]}}{\binom{N}{[n+1]}} \\ &= \mathbb{E}[X_{q,M_q}(n+1)] \end{aligned}$$

so the relation is solved. \square

We can now consider the remaining cluster sizes. For $c = 1, \dots, M_q - 1$, we have;

$$X_{q,c}(n+1) = X_{q,c}(n) + \begin{cases} -1, & \text{if we sample from a cluster of size } c \\ 0, & \text{if we sample from a cluster not of size } c \text{ or } c+1 \\ 1, & \text{if we sample from a cluster of size } c+1 \end{cases}$$

or equivalently

$$X_{q,c}(n+1) = X_{q,c}(n) + \begin{cases} -1, & \text{with probability: } \frac{cX_{q,c}(n)}{N-n} \\ 0, & \text{with probability: } \frac{N-n-cX_{q,c}(n)-(c+1)X_{q,c+1}(n)}{N-n} \\ 1, & \text{with probability: } \frac{(c+1)X_{q,c+1}(n)}{N-n} \end{cases}$$

Taking expectations as before gives us the recurrence relation;

$$\mathbb{E}(X_{q,c}(n+1)) = \frac{(N-n-c)\mathbb{E}(X_{q,c}(n))}{N-n} + \frac{(c+1)\mathbb{E}(X_{q,c+1}(n))}{N-n} \quad (17)$$

This is an inhomogeneous recurrence relation, however the homogeneous component given by;

$$\mathbb{E}(X_{q,c}(n+1)) = \frac{(N-n-c)\mathbb{E}(X_{q,c}(n))}{N-n}$$

which is solved by:

$$\mathbb{E}[X_{q,c}(n)] = X_{q,c}(0) \frac{\binom{N-c}{n}}{\binom{N}{n}}$$

However, because theses are inhomogeneous, we must consider the full system given by:

$$\begin{bmatrix} \mathbb{E}(X_{q,1}(n+1)) \\ \mathbb{E}(X_{q,2}(n+1)) \\ \vdots \\ \vdots \\ \mathbb{E}(X_{q,M_q-1}(n+1)) \\ \mathbb{E}(X_{q,M_q}(n+1)) \end{bmatrix} = \begin{bmatrix} \frac{N-n-1}{N-n} & \frac{2}{N-n} & \cdot & \cdot & \cdot & 0 & 0 \\ 0 & \frac{N-n-2}{N-n} & \cdot & \cdot & \cdot & 0 & 0 \\ \cdot & 0 & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & 0 & 0 & \cdot & \cdot & \frac{N-n-(M_q-1)}{N-n} & \frac{M_q}{N-n} \\ 0 & 0 & 0 & \cdot & \cdot & 0 & \frac{N-n-M_q}{N-n} \end{bmatrix} \begin{bmatrix} \mathbb{E}(X_{q,1}(n)) \\ \mathbb{E}(X_{q,2}(n)) \\ \vdots \\ \vdots \\ \mathbb{E}(X_{q,M_q-1}(n)) \\ \mathbb{E}(X_{q,M_q}(n)) \end{bmatrix} \quad (18)$$

with initial conditions given by;

$$(\mathbb{E}(X_{q,1}(0)), \mathbb{E}(X_{q,2}(0)), \dots, \mathbb{E}(X_{q,M_q-1}(0)), \mathbb{E}(X_{q,M_q}(0)))^T = \left(\alpha_{q,1} N, \frac{\alpha_{q,2} N}{2}, \dots, \frac{\alpha_{q,M_q-1} N}{M-1}, \frac{\alpha_{q,M_q} N}{M} \right)^T \quad (19)$$

Theorem 8.2. *The above system of recurrence relations equipped with initial conditions () is solved by:*

$$\mathbb{E}[X_{q,c}(n)] = \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-k+c}}{\binom{N}{n}}$$

Proof. We first note that when $n = 0$ our proposed solution reduces to:

$$\mathbb{E}[X_{q,c}(0)] = X_{q,k}(0) \frac{\binom{c}{0} \binom{N-c}{0}}{\binom{N}{0}} + \sum_{k=c+1}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{c-k}}{\binom{N}{0}} = X_{q,k}(0)$$

where terms for $k > c$ equal 0 by definition of the binomial coefficient for negative numbers. Having shown this is true for $c = M_q$, we assume this is true for some a general $c = 2, \dots, M_q$ and show that the resulting recurrence relation for $c - 1$ is also solved by this solution.

The recurrence relation for $c - 1$ is given by;

$$\mathbb{E}(X_{q,c-1}(n+1)) = \frac{(N-n-(c-1))\mathbb{E}(X_{q,c-1}(n))}{N-n} + \frac{c \cdot \mathbb{E}(X_{q,c}(n))}{N-n}$$

Using the inductive hypothesis, we substitute our result for $X_{q,c}(n)$ and find;

$$\mathbb{E}(X_{q,c-1}(n+1)) = \frac{(N-n-(c-1))\mathbb{E}(X_{q,c-1}(n))}{N-n} + \frac{c}{N-n} \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-k+c}}{\binom{N}{n}}$$

We then aim to show that this relation is solved by;

$$\mathbb{E}(X_{q,c-1}(n)) = \sum_{k=c-1}^{M_q} X_{q,k}(0) \frac{\binom{k}{c-1} \binom{N-k}{n-k+c-1}}{\binom{N}{n}}$$

We have already seen that this proposed solution satisfies the initial condition for the above so need only show that it satisfies the relation. We substitute (x) into the relation and find the right hand

side as:

$$\begin{aligned}
& \frac{(N-n-(c-1))}{N-n} \sum_{k=c-1}^{M_q} X_{q,k}(0) \frac{\binom{k}{c-1} \binom{N-k}{n-k+c-1}}{\binom{N}{n}} + \frac{c}{N-n} \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c} \binom{N-k}{n-k+c}}{\binom{N}{n}} \\
&= \frac{N-n-c+1}{N-n} \cdot X_{q,c-1}(0) \frac{\binom{N-c+1}{n}}{\binom{N}{n}} + \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{N-n-c+1}{N-n} \cdot \frac{\binom{k}{c-1} \binom{N-k}{n-k+c-1}}{\binom{N}{n}} + \frac{c}{N-n} \cdot \frac{\binom{k}{c} \binom{N-k}{n-k+c}}{\binom{N}{n}} \right] \\
&= X_{q,c-1}(0) \frac{\binom{N-c+1}{[n+1]}}{\binom{N}{[n+1]}} + \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{\binom{k}{c-1} \binom{N-k}{[n+1]-k+c-1}}{\binom{N}{[n+1]}} \frac{n-k+c}{n+1} + \frac{\binom{k}{c-1} \binom{N-k}{n-k+c}}{\binom{N}{[n+1]}} \frac{k-c+1}{n+1} \right] \\
&= X_{q,c-1}(0) \frac{\binom{N-c+1}{[n+1]}}{\binom{N}{[n+1]}} + \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c-1}}{\binom{N}{[n+1]}} \left[\frac{\binom{N-k}{[n+1]-k+c-1} (n-k+c) + \binom{N-k}{[n+1]-k+c-1} (k-c+1)}{n+1} \right] \\
&= X_{q,c-1}(0) \frac{\binom{N-c+1}{[n+1]}}{\binom{N}{[n+1]}} + \sum_{k=c}^{M_q} X_{q,k}(0) \frac{\binom{k}{c-1} \binom{N-k}{[n+1]-k+c-1}}{\binom{N}{[n+1]}} \\
&= \sum_{k=c-1}^{M_q} X_{q,k}(0) \frac{\binom{k}{c-1} \binom{N-k}{[n+1]-k+c-1}}{\binom{N}{[n+1]}} \\
&= \mathbb{E}[X_{q,c-1}(n+1)]
\end{aligned}$$

so given the result is true for c , the proposed result satisfies the resulting recurrence relation for $c-1$. In conjunction with the proof for $c = M_q$, the result holds for the full system. \square

8.2 Population Variance and Covariance

We have:

$$\begin{aligned}
X_{q,c}(n)^2 &= \left(\sum_{k=c}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{k-c}(n) \right)^2 \\
&= \sum_{k=c}^{M_q} [X_{q,k}(0) \mathbb{1}_{i_k,k,q}^{k-c}(n) + X_{q,k}(0)(X_{q,k}(0)-1) \mathbb{1}_{i_k,k,q}^{k-c}(n) \mathbb{1}_{j_k,k,q}^{k-c}(n)] \\
&\quad + \sum_{k=c}^{M_q-1} \sum_{l=k+1}^{M_q} X_{q,k}(0) X_{q,l}(0) \mathbb{1}_{i_k,k,q}^{k-c}(n) \mathbb{1}_{i_l,l,q}^{l-c}(n)
\end{aligned}$$

So taking expectations we find:

$$\begin{aligned}
\mathbb{E}[X_{q,c}(n)^2] &= \sum_{k=c}^{M_q} X_{q,k}(0) \left[\frac{\binom{k}{k-c} \binom{N-k}{n-k+c}}{\binom{N}{n}} - \frac{\binom{k}{k-c} \binom{k}{k-c} \binom{N-2k}{n-2k+2c}}{\binom{N}{n}} \right] \\
&\quad + \sum_{k=c}^{M_q} \sum_{l=c}^{M_q} X_{q,k}(0) X_{q,l}(0) \frac{\binom{k}{k-c} \binom{l}{l-c} \binom{N-k-l}{n-k-l+2c}}{\binom{N}{n}}
\end{aligned}$$

so using the binomial coefficient symmetries discussed in (here), we have:

$$\mathbb{E}[X_{q,c}(n)^2] = \mathbb{E}[S_{q,c}(N-n)^2]$$

and hence, as:

$$\mathbb{E}[X_{q,c}(n)]^2 = \mathbb{E}[S_{q,c}(N-n)]^2$$

the result follows.

For the joint expectation we have:

$$\begin{aligned}
X_{q,c_1}(n)X_{q,c_2}(n) &= \left(\sum_{k=c_1}^{M_q} \sum_{i_k=1}^{X_{q,k}(0)} \mathbb{1}_{i_k,k,q}^{k-c_1}(n) \right) \left(\sum_{l=c_2}^{M_q} \sum_{i_l=1}^{X_{q,l}(0)} \mathbb{1}_{i_l,l,q}^{l-c_2}(n) \right) \\
&= \sum_{k=c_2}^{M_q} X_{q,k}(0) \mathbb{1}_{i_k,k,q}^{k-c_1}(n) \mathbb{1}_{i_k,k,q}^{k-c_2}(n) + X_{q,k}(0)(X_{q,k}(0)-1) \mathbb{1}_{i_k,k,q}^{k-c_1}(n) \mathbb{1}_{j_k,k,q}^{k-c_2}(n) \\
&\quad + \sum_{k \neq l} X_{q,k}(0)X_{q,l}(0) \mathbb{1}_{i_k,k,q}^{k-c_1}(n) \mathbb{1}_{i_l,l,q}^{k-c_2}
\end{aligned}$$

Then after taking expectations we find:

$$\begin{aligned}
\mathbb{E}[X_{q,c_1}(n)X_{q,c_2}(n)] &= \sum_{k=c_1}^{M_q} \sum_{l=c_2}^{M_q} X_{q,k}(0)X_{q,l}(0) \frac{\binom{k}{k-c_1} \binom{l}{l-c_2} \binom{N-k-l}{n-k-l+c_1+c_2}}{\binom{N}{n}} \\
&\quad - \sum_{k=c_2}^{M_q} X_{q,k}(0) \frac{\binom{k}{k-c_1} \binom{k}{k-c_2} \binom{N-2k}{n-2k+c_1+c_2}}{\binom{N}{n}} \\
&= \mathbb{E}[S_{q,c_1}(N-n)S_{q,c_2}(N-n)]
\end{aligned}$$

so using that:

$$\mathbb{E}[X_{q,c_1}(n)]\mathbb{E}[X_{q,c_2}(n)] = \mathbb{E}[S_{q,c_1}(N-n)]\mathbb{E}[S_{q,c_2}(N-n)]$$

as well as inferring that the negative sum disappears for $q_1 \neq q_2$ we reach the result.

8.3 Variance of Species Estimators

We have found:

$$\hat{S}_n = \sum_{i=1}^S \mathbb{1}_{\{X_i>0\}}(n)$$

so that:

$$\begin{aligned}
\hat{S}_n^2 &= \sum_{i=1}^S \mathbb{1}_{\{X_i>0\}}(n)^2 + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \mathbb{1}_{\{X_i>0\}}(n) \mathbb{1}_{\{X_j>0\}}(n) \\
&= \sum_{i=1}^S \mathbb{1}_{\{X_i>0\}}(n) + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \mathbb{1}_{\{X_i,X_j>0\}}(n)
\end{aligned}$$

So that in the case of sampling without replacement,

$$\mathbb{E}[\hat{S}_n^2] = \sum_{i=1}^S \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right) + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} - \frac{\binom{N-N_j}{n}}{\binom{N}{n}} + \frac{\binom{N-N_i-N_j}{n}}{\binom{N}{n}} \right)$$

Then we calculate:

$$\begin{aligned}
\mathbb{E}[\hat{S}_n]^2 &= \left(\sum_{i=1}^S \mathbb{P}(X_i > 0) \right)^2 \\
&= \sum_{i=1}^S \mathbb{P}(X_i > 0)^2 + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \mathbb{P}(X_i > 0) \mathbb{P}(X_j > 0) \\
&= \sum_{i=1}^S \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right)^2 + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right) \left(1 - \frac{\binom{N-N_j}{n}}{\binom{N}{n}} \right)
\end{aligned}$$

and therefore, we find:

$$\begin{aligned}
\text{Var}(\hat{S}_n) &= \mathbb{E}[\hat{S}_n]^2 - \mathbb{E}[\hat{S}_n]^2 \\
&= \sum_{i=1}^S \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \left(1 - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \right) + 2 \sum_{i=1}^{S-1} \sum_{j=i+1}^S \left(\frac{\binom{N-N_i-N_j}{n}}{\binom{N}{n}} - \frac{\binom{N-N_i}{n}}{\binom{N}{n}} \frac{\binom{N-N_j}{n}}{\binom{N}{n}} \right) \\
&= \sum_{i=1}^S \left(\frac{\binom{N-N_i}{n}}{\binom{N}{n}} - \frac{\binom{N-2N_i}{n}}{\binom{N}{n}} \right) + \sum_{i=1}^S \sum_{j=1}^S \left(\frac{\binom{N-N_i-N_j}{n}}{\binom{N}{n}} - \frac{\binom{N-N_i}{n} \binom{N-N_j}{n}}{\binom{N}{n}^2} \right)
\end{aligned}$$

The results for the with replacement case follow by noting that the probabilities are then given by:

$$\mathbb{P}(X_i > 0) = 1 - \left(1 - \frac{N_i}{N} \right), \text{ and } \mathbb{P}(X_i > 0, X_j > 0) = 1 - \left(1 - \frac{N_i}{N} \right) \left(1 - \frac{N_j}{N} \right) + \left(1 - \frac{N_i+N_j}{N} \right)$$

8.4 Harmonic Inequality

Geometrically, by considering the Riemann sums for the harmonic series, integral bounds are given as:

$$1 + \frac{1}{2} + \dots + \frac{1}{n} > \int_1^{n+1} \frac{1}{x} dx$$

and:

$$\frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n} < \int_1^n \frac{1}{x} dx$$

Therefore, for

$$H_m^n = H_n - H_{m-1} = \frac{1}{m} + \frac{1}{m+1} + \dots + \frac{1}{n} \quad (20)$$

we have:

$$\int_m^{n+1} \frac{1}{x} dx < H_m^n < \int_{m-1}^n \frac{1}{x} dx$$

which gives the proposed bounds.