# HOW MANY IID SAMPLES
# DOES IT TAKE TO SEE ALL
# THE BALLS IN A BOX?

by

Thomas M. Sellke
Purdue University

Department of Statistics
Purdue University

◇

# How Many IID Samples Does it Take to See All the Balls in a Box?

**Thomas M. Sellke**

**Purdue University**

## Abstract

Suppose a box contains $m$ balls, numbered from 1 to $m$. A random number of balls are drawn from the box, their numbers are noted, and the balls are then returned to the box. This is done repeatedly, with the sample sizes being iid. Let $X$ be the number of samples needed to see all the balls. This paper uses Markov-chain coupling to derive a simple but typically very accurate approximation for $EX$ in terms of the sample size distribution. The approximation formula generalizes the formula found by Polya (1930) for the special case of fixed sample sizes.

## 1. Introduction

Suppose we have a box containing $m$ identical white balls. Let $K_1, K_2, \ldots$ be iid random variables taking positive integer values. We randomly sample $K_1 \wedge m$ balls without replacement, paint the sampled balls red, and return them to the box. Then $K_2 \wedge m$ balls are sampled, the white ones are painted red, and all are returned to the box, etc. Let $X$ be the number of samples needed to paint all the balls red. When $\max\limits_{1 \leq i \leq X} K_i$ is with high probability smaller than $m/2$, say, then a good approximation for $EX$ is given by

$$
(1.1) \qquad \frac{\sum\limits_{i=0}^{m-1} \frac{1}{m-i}}{\sum\limits_{i=0}^{m-1} \frac{1}{m-i} P\{K > i\}} + \frac{\sum\limits_{r=1}^{m-1} \frac{1}{m-r} P\{K > r\} \sum\limits_{j=1}^{r} \frac{1}{m-j+1}}{\left[ \sum\limits_{i=0}^{m-1} \frac{1}{m-i} P\{K > i\} \right]^2}.
$$

(A $K$ without a subscript represents a generic $K_i$.) For instance, if $m = 10$ and the $(K_i - 1)$'s are binomial $(4, \frac{1}{2})$, then the true value of $EX$ is 8.8937, while (1.1) gives 8.8933. For $m = 20$ and $(K_i - 1)$'s which are binomial $(4, \frac{1}{2})$, the values are $EX = 22.90753529760067$ and $(1.1) = 22.90753529760074$. For $m = 10$ and $K_i$'s which are uniformly distributed on $\{1, 2, 3, 4, 5\}$, $EX = 8.74239$ and $(1.1) = 8.74236$. For $m = 20$, the values are $EX = 22.740208948996$ and $(1.1) = 22.740208948981$. (The true values

were computed by Jacek Dmochowski using exact recursive formulas suggested by Larry Shepp.) It is difficult to determine the size of the approximation error when $m$ is much larger than 20. Even with double-precision computation, the round-off error seems to dominate the true approximation error.

The justification of formula (1.1) involves Wald's identity (which gives the first term of (1.1) as a first approximation to $EX$) and coupling (which gives the second term of (1.1) as a correction for "boundary overshoot error" in the first term). A bound on the difference between $EX$ and (1.1) can be given in terms of, say, $P\{\max_{i \leq X} K_i > \frac{m}{2}\}$ and the probability that a certain Markov-chain coupling is unsuccessful.

Section 4 gives some explicit bounds on the approximation error. These bounds are generally very crude, but they show that the approximation error converges to zero faster than $\exp(-m^{\frac{1}{3}})$ as $m \to \infty$ when the (fixed) $K$-distribution has a finite moment generating function near 0. If the $K_i$'s are bounded, or if the hazard function of the $K_i$'s is bounded below by $\delta > 0$, then the approximation error converges to zero exponentially (in $m$) as $m \to \infty$.

Section 5 presents a generalization of (1.1) applicable to the case where some of the balls in the box are red to begin with, and where only a specified number of the white balls need to be painted red.

Polya (1930) found an approximation formula for $EX$ when the sample size $K$ is constant. Formula (1.1) agrees with Polya's formula in this special case. Polya (1930) indicated how to prove that the approximation error converges to zero as $m \to \infty$, but with no rate. However, he also showed that the approximation error is exactly $(-1)^m m(m-1)/\{(2m-1)^2 \binom{2m-2}{m}\}$ when $K \equiv 2$, so he probably suspected that his approximation would be extremely good in general for large $m$. His methods were completely different from ours.

## 2. Application of the Wald Identity

Assume the balls are numbered $1, 2, \ldots, m$. If balls are sampled one at a time, with replacement, then analyzing the number $\tau$ of single-ball draws needed to see all the balls is easy: it's just the standard coupon collector problem, as described for example by Ross (1988) or Feller (1968). Once $j$ balls have been seen, the number of additional draws

2

needed to see a new ball is a geometric $(\frac{m-j}{m})$ random variable independent of everything that has come before. (Here, $X \sim$ geometric $(p)$ means $P\{X = k\} = (1-p)^{k-1}p$ for $k = 1, 2, \ldots$, with $EX = p^{-1}$). Adding expectations yields

(2.1)
$$E\tau = \sum_{j=0}^{m-1} \frac{m}{m-j}.$$

(Note that (1.1) = (2.1) when $P\{K = 1\} = 1$.)

More generally, if $K$ is a random variable (independent of the drawing process) taking positive integer values, the number $D$ of single-ball draws needed to see $K \wedge m$ distinct balls has expectation

(2.2)
$$ED = \sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\},$$

since $P\{K > i\}$ is the probability that the geometric $(\frac{m-i}{m})$ number of draws needed to get the $(i+1)^{st}$ distinct ball is included in $D$.

Suppose our successive samples of balls are obtained as follows. First we see the value of $K_1$. Then we draw balls, one at a time, *with replacement*, until $K_1 \wedge m$ distinct balls have been obtained. Let $D_1$ be the number of draws needed to obtain the required $K_1 \wedge m$ distinct balls. Then we see $K_2$ and draw balls one at a time with replacement, until $K_2 \wedge m$ distinct balls have been obtained to form the second sample, etc. Let $D_i$ be the number of draws needed to obtain the $K_i \wedge m$ distinct balls in the $i^{th}$ sample. Let $\mathcal{F}_i$ be the $\sigma$-field generated by all observations made while generating the first $i$ samples. Thus, $\mathcal{F}_i$ is generated by $K_1, K_2, \ldots, K_i$ *and* by the sequence of $D_1 + \ldots + D_i$ draws needed to obtain the first $i$ samples. The $K_i$'s are assumed to be iid. The numbers on the balls obtained on successive draws are iid, uniformly distributed on $\{1, 2, \ldots, m\}$, and independent of the $K_i$'s. The $D_i$'s are also iid, and $K_{i+1}$ and $D_{i+1}$ are independent of $\mathcal{F}_i$.

Again, let $X$ be the number of samples needed to see all the balls. The event $\{X \leq i\}$ is $\mathcal{F}_i$ measurable, since we know after the first $D_1 + \ldots + D_i$ draws whether or not we have seen all the balls. Thus, $X$ is an $\{\mathcal{F}_i\}$ stopping time. By Wald's identity/equality/lemma/relation and the fact that the $D_i$'s are iid,

(2.3)
$$E(D_1 + \ldots + D_X) = EX \, ED_1.$$

3

With $\tau$ as above, note that

$$X = \inf\{i : D_1 + \ldots + D_i \geq \tau\}.$$

Thus, the sum $\sum_1^X D_i$ and $\tau$ are equal except for the "overshoot" given by the number of draws needed to complete the last sample after the last new ball has been obtained on the $\tau^{th}$ draw. Let $V$ be the size of this overshoot, so that

(2.4)
$$V = \left(\sum_1^X D_i\right) - \tau.$$

From (2.1), (2.2), (2.3), and (2.4) we get

(2.5)
$$EX = \frac{E\sum_1^X D_i}{ED_1} = \frac{E\tau + EV}{ED_1} = \frac{\sum_{i=0}^{m-1} \frac{m}{m-i} + EV}{\sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\}}.$$

Let $A_\infty$ be the number of distinct balls already in the last sample after the $\tau^{th}$ draw. The same "coupon collector" argument used for (2.2) shows that

$$E(V|A_\infty = j) = \sum_{r=j}^{m-1} \frac{m}{m-r} P\{K_X > r|A_\infty = j\}.$$

But the conditional distribution of $K_X$, given that $A_\infty = j$, is exactly the same as that of $K_1$ given that $K_1 \geq j$. Thus,

(2.6)
$$E(V|A_\infty = j) = \sum_{r=j}^{m-1} \frac{m}{m-r} P\{K > r|K \geq j\}.$$

(We interpret (2.6) as zero if $j = m$.) If we can get a good approximation for distribution of $A_\infty$, we can combine this with (2.6) to approximate $EV$ in (2.5). The next section will show that

(2.7)
$$P\{A_\infty = j\} \approx \frac{\frac{m}{m-j+1} P\{K \geq j\}}{\sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\}}.$$

4

Combining (2.7) with (2.6) yields

$$(2.8) \qquad EV = \sum_{j=1}^{m-1} E(V|A_\infty = j)P\{A_\infty = j\}$$

$$\approx \frac{\displaystyle\sum_{j=1}^{m-1}\sum_{r=j}^{m-1} \frac{m}{m-r}P\{K > r|K \geq j\}\frac{m}{m-j+1}P\{K \geq j\}}{\displaystyle\sum_{i=0}^{m-1} \frac{m}{m-i}P\{K > i\}}$$

$$= \frac{\displaystyle\sum_{r=1}^{m-1} \frac{m}{m-r}P\{K > r\} \sum_{j=1}^{r} \frac{m}{m-j+1}}{\displaystyle\sum_{i=0}^{m-1} \frac{m}{m-i}P\{K > i\}}.$$

Substituting this approximation for $EV$ into (2.5) and cancelling some factors of $m$ yields (1.1).

## 3. Approximation of $P\{A_\infty = j\}$ by Coupling

This section will explain the approximation (2.7) above:

$$(3.1) \qquad P\{A_\infty = j\} \approx \frac{\frac{m}{m-j+1}P\{K_1 \geq j\}}{\displaystyle\sum_{i=0}^{m-1} \frac{m}{m-i}P\{K > i\}}, \quad j = 1, 2, \ldots, m.$$

The idea is to construct a finite-state-space Markov chain with $m$ absorbing states labelled $1, 2, \ldots, m$ for which $A_\infty$ equals the label of the state where absorption occurs. This chain is coupled to an approximating chain for which the distribution of the absorbing state is given (modulo truncation) by the right side of (3.1). (For a general account of the coupling technique, see Lindvall (1992).)

For the purposes of this section, it will be better to use a recipe for generating the required samples which is different from that of the previous section. To start, draw a single ball from the box. Then ask whether the next ball should continue the first sample. The probability of continuing the first sample should be

$$c_1 =: \frac{P\{K \wedge m > 1\}}{P\{K \wedge m \geq 1\}}. \qquad (3.2)$$

If we decide to continue the first sample, draw another ball from the box without replacing the first ball. With probability $h_1 =: 1 - c_1$, we decide to end the first sample with the

5

first ball. In this case, we return the first ball to the box and then draw the first ball of the second sample.

The general rule for sampling goes as follows. If the number of balls already in the current sample is $a$, we decide to continue this sample with the next ball with probability

$$c_a =: \frac{P\{K \wedge m > a\}}{P\{K \wedge m \geq a\}}. \tag{3.3}$$

With probability $h_a =: 1 - c_a$, we end the current sample, return all balls in this current sample to the box, and then draw the first ball of the next sample. Balls are returned to the box only when we decide that the current sample is complete and that it's time to start a new one. It should be obvious that this protocol really does generate independent samples whose common sample size distribution is as required.
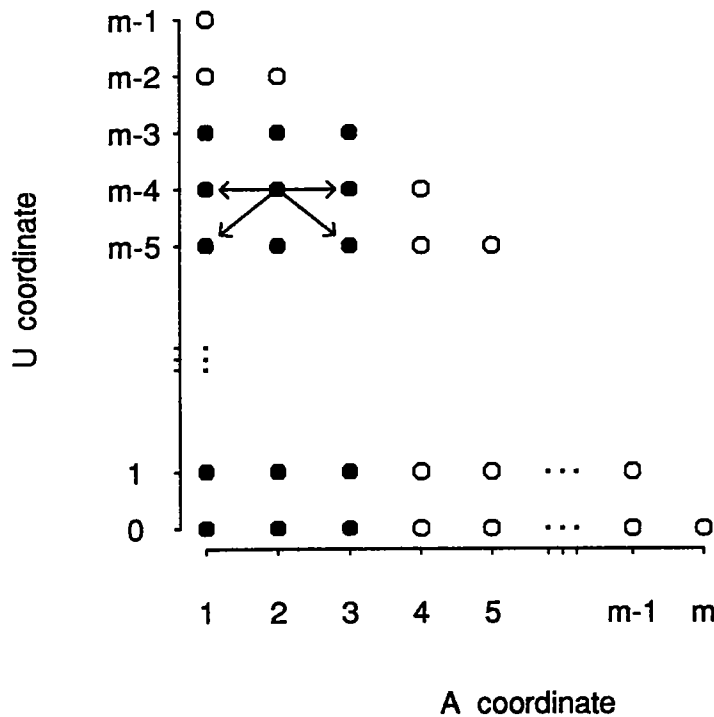


Figure 1. The dots, both solid and hollow, show the state space of the $(A_n, U_n)$ Markov chain. For $b = 3$, the solid dots show the state space of the $(A_n^{(b)}, U_n^{(b)})$ chain, which starts on $U$-level $m - b$. The arrows show the possible transitions out of state $(A, U) = (2, m-4)$. The states on the bottom row are absorbing.

6

Let $\tilde{A}_n$ be the number of balls *already* in the current sample after the $n^{th}$ draw. Let $\tilde{U}_n$ be the number of previously *unsampled* (i.e., white) balls left after the $n^{th}$ draw. Let $T_0$ be the number of draws needed to get all balls at least once (according to the sampling protocol of this section), so that $T_0 = \inf\{n : \tilde{U}_n = 0\}$. Define

$$(3.4) \qquad\qquad A_n =: \tilde{A}_{n \wedge T_0} \text{ and } U_n =: \tilde{U}_{n \wedge T_0}.$$

Then $\{(A_n, U_n)\}_{n=1}^{\infty}$ is a Markov chain with state space

$$S_m =: \{(a, u) : a \in \{1, 2, \ldots, m\}, \ u \in \{0, 1, \ldots, m-1\}, \ a + u \le m\}.$$

Since we start with $n = 1$, the starting state is $(1, m-1)$. For $u \ge 1$, the transition probabilities are

$$(3.5) \qquad P\{(A_{n+1}, U_{n+1}) = (a', u') | (A_n, U_n) = (a, u)\}$$

$$\begin{aligned}
&= h_a \frac{m-u}{m} && \text{for } (a', u') = (1, u) \\
&= h_a \frac{u}{m} && \text{for } (a', u') = (1, u-1) \\
&= c_a \frac{m-a-u}{m-a} && \text{for } (a', u') = (a+1, u) \\
&= c_a \frac{u}{m-a} && \text{for } (a', u') = (a+1, u-1) \\
&= 0 && \text{otherwise.}
\end{aligned}$$

Since the $(A_n, U_n)$ chain stops at time $T_0$, states of the form $(a, 0)$ are absorbing. Note that $A_{T_0} = A_{\infty}$, where $A_{\infty}$ is as defined in Section 2.

Now we define the approximating Markov chain. Fix $b \in \{1, 2, \ldots, m-1\}$. Let $(A_n^{(b)}, U_n^{(b)})$ be a Markov chain on $S_m$ whose transition probabilities are as above, for the $(A_n, U_n)$ chain, but with $K \wedge b$ in place of $K$. Thus,

$$c_a^{(b)} = \begin{cases} c_a & \text{if } a < b \\ 0 & \text{if } a \ge b \end{cases}$$

and $h_a^{(b)} = 1 - c_a^{(b)}$. Note that the two chains $(A_n, U_n)$ and $(A_n^{(b)}, U_n^{(b)})$ have the same transition probabilities out of states $(a, u)$ for which $a < b$. Our coupling construction below will work well when we can choose a value of $b$ not too close to $m$ but with $P\{K > b\}$ small enough to make $P\{\max_{i \le X} K_i > b\}$ negligible.

The $\{(A_n^{(b)}, U_n^{(b)})\}_{n=1}^\infty$ chain will be started with $U_1^{(b)} = m - b$. The value of $A_1^{(b)}$ will be random, with

$$(3.6) \qquad P\{A_1^{(b)} = a\} = \frac{\frac{m}{m-a+1}P\{K \geq a\}}{\sum\limits_{i=0}^{b-1} \frac{m}{m-i}P\{K > i\}}, \quad 1 \leq a \leq b.$$

(Compare (3.6) to (3.1). The distribution for $A_1^{(b)}$ in (3.6) is the distribution on the right side of (3.1), conditioned to be $\leq b$.)

For $u \in \{0, 1, \dots, m - b\}$, define

$$(3.7) \qquad T_u^{(b)} = \inf\{n : U_n^{(b)} = u\}.$$

Thus, $T_u^{(b)}$ is the time at which the $(A_n^{(b)}, U_n^{(b)})$ chain drops down to $U$-level $u$. In analogy with previous notation, set $A_\infty^{(b)} = A_{T_0^{(b)}}^{(b)}$, so that $(A_\infty^{(b)}, 0)$ is the state where the $(A_n^{(b)}, U_n^{(b)})$ chain is absorbed.

Here is the key result for understanding formula (3.1):

**Proposition 3.1.** If $U_1^{(b)} = m - b$ and $A_1^{(b)}$ has distribution (3.6), then $A_{T_u^{(b)}}^{(b)}$ has distribution (3.6) for *all* $u \in \{0, 1, \dots, m - b\}$. In particular, $A_\infty^{(b)}$ has distribution (3.6).

**Proof.** Define $\{B_n^{(b)}\}_{n=1}^\infty$ to be a Markov chain on state space $\{1, 2, \dots, b\}$ which acts like a non-stopped version of $A_n^{(b)}$. Thus, the transition probabilities for the $B_n^{(b)}$ chain are

$$(3.8) \qquad P\{B_{n+1}^{(b)} = a' | B_n^{(b)} = a\}$$

$$= c_a^{(b)} = \frac{P\{K \wedge b > a\}}{P\{K \wedge b \geq a\}} \quad \text{if } a' = a + 1$$
$$= h_a^{(b)} = \frac{P\{K \wedge b = a\}}{P\{K \wedge b \geq a\}} \quad \text{if } a' = 1$$
$$= 0 \qquad \text{otherwise.}$$

It is well-known (and trivial to check) that the stationary distribution of this "renewal age process" chain is

$$(3.9) \qquad \pi_a = \frac{P\{K \geq a\}}{\sum\limits_{i=1}^{b} P\{K \geq i\}} = \frac{P\{K \geq a\}}{E(K \wedge b)}, \quad 1 \leq a \leq b.$$

8

Define another Markov chain $(B_n^{(b)}, W_n^{(b)})$ by having $B_n^{(b)}$ as above and the $W_n^{(b)}$'s a sequence of Bernoulli random variables. Set $W_1^{(b)} \equiv 1$. Conditional on the entire path

$$\underline{B}^{(b)} =: (B_1^{(b)}, B_2^{(b)}, \dots)$$

of the $B_n^{(b)}$ chain, the $W_n^{(b)}$'s, $n \geq 2$, are to be independent, with

$$(3.10) \qquad P\{W_n^{(b)} = 1 | \underline{B}^{(b)}\} = 1 - P\{W_n^{(b)} = 0 | \underline{B}^{(b)}\} = \frac{d}{m - B_n^{(b)} + 1}$$

for some constant $d$, $0 < d \leq m - b$.

Here is how to think about the $(B_n^{(b)}, W_n^{(b)})$ chain in terms of Figure 1. Suppose the $(A_n^{(b)}, U_n^{(b)})$ chain is stuck on some $U$-level $d$, in the sense that it is artificially held at $U$-level $d$ whenever it "wants" to drop down to $U$-level $d - 1$. If $B_n^{(b)}$ is the $A$ coordinate of this process and $W_n^{(b)}$ is an indicator of the event that the process "wanted" to drop its $U$-level from $d$ to $d - 1$ between times $n - 1$ and $n$, then $(B_n^{(b)}, W_n^{(b)})$ is as above.

Now consider the "embedded chain" whose consecutive states are the consecutive states of the $(B_n^{(b)}, W_n^{(b)})$ chain at which $W_n^{(b)} = 1$. (The times at which $W_n^{(b)} = 0$ are skipped.) The stationary distribution for $B_n^{(b)}$ in this embedded chain is

$$(3.11) \qquad \pi_a^{\text{emb}} = \frac{\frac{1}{m-a+1}\pi_a}{\sum\limits_{i=1}^{u} \frac{1}{m-i+1}\pi_i} = \frac{\frac{1}{m-a+1}P\{K \geq a\}}{\sum\limits_{i=0}^{u} \frac{1}{m-i}P\{K > i\}}, \quad 1 \leq a \leq b.$$

(Note that (3.11) and (3.6) are the same distribution.) The reasoning for (3.11) is as follows. The stationary distribution of an ergodic Markov chain gives the long-run fraction of the time that the chain spends in each state. By (3.9), (3.10), and the strong law of large numbers, the long-run fraction of the time that the $(B_n^{(b)}, W_n^{(b)})$ chain spends in state $(a, 1)$ equals $\pi_a \frac{d}{(m-a+1)}$. Thus, the long-run fraction of the time that the *embedded* chain spends in a state $(a, 1)$ is given by (3.11).

Recall that if the starting state of an ergodic Markov chain is chosen according to the stationary distribution, then the state one time unit later will also have the stationary distribution.

Now compare the $(A_n^{(b)}, U_n^{(b)})$ chain with the $(B_n^{(b)}, W_n^{(b)})$ chain. The differences $U_{n-1}^{(b)} - U_n^{(b)}$ are Bernoulli random variables. For $n \leq T_{m-b-1}^{(b)}$, (i.e., until the $(A_n^{(b)}, U_n^{(b)})$ chain

drops from its starting level $m - b$ to the next lower level),

$$P\{U_{n-1}^{(b)} - U_n^{(b)} = 1 | A_n^{(b)}, \{(A_\ell^{(b)}, U_\ell^{(b)})\}_{\ell=1}^{n-1} = \frac{m - b}{m - A_n + 1},$$

which looks like (3.10) with $d = m - b$. Thus, up until time

$$T_{m-b-1}^{(b)} = \inf\{n : U_{n-1}^{(b)} - U_n^{(b)} = 1\},$$

$(A_n^{(b)}, U_{n-1}^{(b)} - U_n^{(b)})$ is a chain with the same state space and transition probabilities as $(B_n^{(b)}, W_n^{(b)})$. Since $A_{T_{m-b-1}^{(b)}}^{(b)}$ corresponds to the second value of $B_n^{(b)}$ in the embedded chain, $A_{T_{m-b-1}^{(b)}}^{(b)}$ will have the stationary distribution (3.6) = (3.11) if $A_1^{(b)}$ starts out in this stationary distribution.

It follows in the same way that $A_{T_i^{(b)}}^{(b)}$ has distribution (3.6) if $A_{T_{i-1}^{(b)}}^{(b)}$ has distribution (3.6). Thus, Proposition 3.1 follows by induction.

$\square$

It remains to show that the distribution of $A_\infty$ is close to the distribution (3.6) of $A_\infty^{(b)}$ for properly chosen $b$. Although analytic arguments may be possible, the easiest way to achieve this would seem to be to couple the $(A_n, U_n)$ and $(A_n^{(b)}, U_n^{(b)})$ chains so that, with high probability, they end up at the same absorbing state. We start by choosing a starting state for the $(A_n^{(b)}, U_n^{(b)})$ chain as described above, with $u_1^{(b)} = m - b$ and $A_1^{(b)} \sim$ (3.6). Then we let the $(A_n, U_n)$ chain (which always starts in state $(1, m - 1)$) run until $U_n = m - b$. Once the chains are on the same $U$-level, we can sequentially choose one or the other chain to take a step, with the goal of course being to get them to inhabit the same state at the same time. Or, we can let both chains move simultaneously, with the transitions for the two chains being dependent if we want. The only requirement is that, each time a chain takes a step, it must do so according to its own transition probabilities; this will guarantee that the path of each individual chain has the right probabilities. (In particular, the distribution of the absorbing state will be correct.) If we can get the chains to meet, we couple them so that they move together. Since the transition probabilities are the same as long as $A_n$ and $A_n^{(b)}$ are $< b$, the chains stay coupled as long as the common $A$ coordinate is $< b$. A coupling can fall apart, however. If the common $A$-value of the two coupled chains reaches $b$, then with probability $c_b$ the chains will be decoupled after

10

the next step. If one chain has its $U$ coordinate decrease before coupling is achieved on a common $U$-level, we just leave this chain alone for a while and run the other chain until its $U$ coordinate also drops, after which we try to achieve a coupling at this next lower $U$-level.

The next section will use a certain coupling strategy, in which we sequentially choose one or the other chain to take a step, to derive explicit error bounds for (3.1) and (1.1). When the discrete hazard function $h_a$ of $K$ is bounded below by a constant $\delta > 0$, another coupling strategy involving simultaneous and dependent steps whenever the two chains are on the same $U$-level sometimes gives better error bounds. The next section will also state but not prove some bounds of this type.

## 4. Bounds on the Approximation Error of (3.1) and (1.1)

Let

$$||P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}|| = \sum_{j=1}^m |P\{A_\infty = j\} - P\{A_\infty^{(b)} = j\}|$$

denote the total variation distance between the distributions of $A_\infty$ and $A_\infty^{(b)}$.

**Proposition 4.1.** If $P\{K > b\} = 0$, then

$$||P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}|| \leq 15e^{-m/b}.$$

**Proof.** The total variation distance is obviously bounded by twice the probability that the $(A_n, U_n)$ and $(A_n^{(b)}, U_n^{(b)})$ are not coupled when they are absorbed. (cf Lindvall (1992), page 12.)

Here is a coupling strategy which has probability less than $(15/2)\exp(-m/b)$ that the chains are not coupled at absorption. Order the (at most) $b$ possible states on each $U$-level of the state space $S_m$ as follows:

$$(2, u) \prec (3, u) \prec \ldots \prec (b, u) \prec (1, u).$$

When both chains are on the same $U$-level, run the chain which is "behind" according to this ordering and leave the other chain alone. If the chains meet, couple them. (Since

11

$A_n > b$ is impossible, the coupling will not be broken later.) If one chain drops below the other, run the upper chain until both are again on the same $U$-level.

If no decrease in a $U$ coordinate occurs for the first $b - 1$ steps taken when the Markov chains are both on $U$-level $u$, then the chains are *guaranteed* to couple on this $U$-level. Each time a chain takes a step starting on $U$-level $u$, the probability that its $U$ coordinate will *not* decrease is at least $\frac{(m-b-u)}{(m-b)}$. (See (3.5).) Thus, the probability of no coupling on $U$-level $u$ is at most

$$1 - (\frac{m - b - u}{m - b})^{b-1}.$$

The probability that coupling never occurs on *any* of the $U$-levels $m - b,\, m - b - 1, \ldots, 2, 1$ is at most

$$(4.1) \qquad \prod_{u=1}^{m-b} \{1 - (\frac{m - b - u}{m - b})^{b-1}\} = \prod_{i=0}^{m-b-1} \{1 - (\frac{i}{m - b})^{b-1}\}$$

$$\leq \exp\{- \sum_{i=0}^{m-b-1} (\frac{i}{m - b})^{b-1}\}, \text{ since } 1 - x \leq e^{-x}.$$

But

$$(4.2) \qquad \sum_{i=0}^{m-b-1} (\frac{i}{m - b})^{b-1} = \sum_{i=0}^{m-b} (\frac{i}{m - b})^{b-1} - 1$$

$$> (m - b) \int_0^1 x^{b-1} dx - 1 = \frac{m - b}{b} - 1 = \frac{m}{b} - 2.$$

Thus, $\exp\{2 - (m/b)\}$ bounds the probability that the chains are not coupled at absorption. Since $2e^2 < 15$, the proposition follows:

$$\square$$

We will see that Proposition 4.1 still holds when $P\{K > b\} > 0$, provided we add $2P\{\max_{i \leq X} K_i > b\}$ onto the right side of the inequality in Proposition 4.1. To bound $P\{\max_{i \leq X} K_i > b\}$, it will help to have a bound on $X$.

If we let $\tau$ be the number of single-ball draws (with balls replaced after each draw) needed to see every ball at least once (as in Section 2), then it is obvious that $P\{X \geq t\} \leq P\{\tau \geq t\}$ for all $t$. (Recall we assume $P\{K \geq 1\} = 1$.)

**Lemma 4.1.** $P\{X \geq m^2\} \leq P\{\tau \geq m^2\} < 4m^{\frac{1}{2}} e^{-m/2}$.

12

**Proof.** As was mentioned just before formula (2.1), $\tau$ is a sum of independent geometric $(\frac{m-j}{m})$ random variables, $j = 0, \ldots, m-1$. Thus $\tau - m$ has factorial moment generating function

$$\phi_{\tau-m}(s) =: E(s^{\tau-m}) = \prod_{j=0}^{m-1} \frac{m-j}{m-js} = \prod_{i=1}^{m} \frac{i}{m - (m-i)s}.$$

Setting $s = 1 + \frac{1}{2m}$, we get

$$E(1 + \frac{1}{2m})^{\tau-m} = \prod_{i=1}^{m} (1 + \frac{m-i}{2mi - m + i})$$

$$< \prod_{i=1}^{m} (1 + \frac{m}{2mi - m} - \frac{i}{2mi}) = \prod_{i=1}^{m} (1 + \frac{1}{2i-1} - \frac{1}{2m})$$

$$< \exp\left\{ \sum_{i=1}^{m} (\frac{1}{2i-1} - \frac{1}{2m}) \right\} = \exp(-\frac{1}{2} + \sum_{i=1}^{m} \frac{1}{2i-1}).$$

But

$$\sum_{i=1}^{m} \frac{1}{2i-1} = 1 + \frac{1}{3} + \frac{1}{5} + \ldots + \frac{1}{2m-1}$$

$$< \frac{4}{3} + \frac{1}{2} \int_{4}^{2m} x^{-1} dx < 1 + \frac{(\log m)}{2}.$$

Thus,

$$E(1 + \frac{1}{2m})^{\tau-m} < (me)^{\frac{1}{2}}.$$

However,

$$(1 + \frac{1}{2m})^{m^2-m} > \exp\{(\frac{1}{2m} - \frac{1}{8m^2})(m^2 - m)\} = \exp(\frac{m}{2} - \frac{5}{8})$$

$$(\text{since } \log (1 + x) > x - \frac{x^2}{2} \text{ for } 0 \leq x \leq 1).$$

Thus, by the Markov inequality,

$$P\{\tau \geq m^2\} < (me)^{\frac{1}{2}} e^{\frac{5}{8} - \frac{m}{2}} < 4m^{\frac{1}{2}} e^{-\frac{m}{2}}.$$

$\square$

**Proposition 4.2.** If $P\{K > b\} < \varepsilon$, then

$$\|P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}\| < 15e^{-m/b} + 2m^2\varepsilon + 8m^{\frac{1}{2}} e^{-m/2}.$$

13

**Proof.** Use the same coupling strategy as in the proof of Proposition 4.1. Expression (4.1) (and therefore $(15/2)e^{-m/b}$) bounds the probability of the event "no (initial) coupling occurs *and* $\max\limits_{i\leq X} K_i \leq b$." But by Lemma 4.1,

$$(4.3) \qquad P\{\max\limits_{i\geq X} K_i > b\} \leq m^2 P\{K > b\} + P\{X > m^2\}$$

$$< m^2\varepsilon + 4m^{\frac{1}{2}}e^{-m/2}.$$

Thus the probability of "no (initial) coupling occurs *or* $\max\limits_{i\leq X} K_i > b$" is bounded by

$$(4.4) \qquad \frac{15}{2}e^{-m/b} + m^2\varepsilon + 4m^{\frac{1}{2}}e^{-m/2}.$$

Note that a coupling, once made, is never broken on $\{\max\limits_{i\leq X} K_i \leq b\}$. Thus, (4.4) bounds the probability that the $(A_n, V_n)$ and $(A_n^{(b)}, V_n^{(b)})$ chains are not coupled at absorption. Multiplying (4.4) by 2 gives the desired bound on the total variation distance.

$\square$

Using the bounds on the total variation distance $||P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}||$ given by Proposition 4.1 and 4.2, it is easy to derive bounds on the difference between (1.1) and $EX$. First a few more lemmas.

**Lemma 4.2.** For any $b \in \{1, 2, \ldots, m-1\}$,

$$|EV - \sum_{j=1}^{m} E(V|A_\infty = j)P\{A_\infty^{(b)} = j\}|$$

$$\leq \frac{1}{2}||P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}|| \max_{j} E(V|A_\infty = j),$$

where $E(V|A_\infty = j)$ is set equal to 0 when $P\{K \geq j\} = 0$.

**Proof.** The left side equals

$$|\sum_{j=1}^{m} E(V|A_\infty = j)[P\{A_\infty = j\} - P\{A_\infty^{(b)} = j\}]|,$$

which is less than or equal to the right side.

$\square$

14

**Lemma 4.3.** For any $b \in \{1, \ldots, m-1\}$

$$\left| \frac{\sum_{r=1}^{m-1} \frac{m}{m-r} P\{K > r\} \sum_{j=1}^{r} \frac{m}{m-j+1}}{\sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\}} - \frac{\sum_{r=1}^{b-1} \frac{m}{m-r} P\{K > r\} \sum_{j=1}^{r} \frac{m}{m-j+1}}{\sum_{i=0}^{b-1} \frac{m}{m-i} P\{K > i\}} \right|$$

$$\leq \frac{\sum_{r=b}^{m-1} \frac{m}{m-r} P\{K > r\} \sum_{j=1}^{r} \frac{m}{m-j+1}}{\sum_{i=0}^{b-1} \frac{m}{m-i} P\{K > i\}}$$

*Remark.* The first term in Lemma 4.3 is the approximation (2.8) for $EV$ that we used to get (1.1). The second term in Lemma 4.3 equals $\sum_{j=1}^{m} E(V|A_\infty = j) P\{A_\infty^{(b)} = j\}$.

**Proof.** $\sum_{j=1}^{r} \frac{m}{m-j+1}$ is obviously increasing in $r$, so the first term between absolute value signs is larger than the second term. The first numerator is greater than the second numerator, and the first denominator is greater than the second denominator. Thus, the difference is less than the difference between numerators divided by the second (smaller) denominator.

$\square$

**Proposition 4.3.** The difference between $EX$ and (1.1) is bounded in absolute value by

$$\min_{0 < b < m} \left\{ \frac{1}{2} \| P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\} \| \max_{j \geq 1} \left[ \sum_{r=j}^{m-1} \frac{m}{m-r} P\{K > r | K \geq j\} \right] \right.$$

$$\left. + \frac{\sum_{r=b}^{m-1} \frac{m}{m-r} P\{K > r\} \sum_{j=1}^{r} \frac{m}{m-j+1}}{\sum_{i=0}^{b-1} \frac{m}{m-i} P\{K > i\}} \right\} \bigg/ \sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\}.$$

**Proof.** This bound follows from (2.5) and Lemmas 4.2 and 4.3.

$\square$

Now let's specialize to get bounds which are not quite such a horrendous mess like the bound in Proposition 4.3.

15

**Proposition 4.4.** If $P\{K > b\} = 0$, the difference between $EX$ and (1.1) is bounded in absolute value by

$$\frac{\frac{15}{2}e^{-m/b}\frac{m(b-1)}{m-b}}{\sum\limits_{i=0}^{b-1}\frac{m}{m-i}P\{K > i\}} < \frac{15m(b-1)e^{-m/b}}{2(m-b)EK}.$$

**Proof.** If $P\{K > b\} = 0$, then

$$\max_{j \geq 1}\sum_{r=j}^{m-1}\frac{m}{m-r}P\{K > r | K \geq j\} \leq \frac{m(b-1)}{m-b}.$$

Apply this and Proposition 4.1 to the bound in Proposition 4.3.

$\square$

*Remark.* If $m = 100$ and each $K_i - 1$ is binomial $(4, \frac{1}{2})$, the second error bound in Proposition 4.4 equals $1.96 \times 10^{-4}$ (with $b = 5$ and $EK = 3$).

**Proposition 4.5.** If $P\{K > k\} \leq Ce^{-\alpha k}$ for $k \leq b$, $C > 0$, and $\alpha > 0$, then the difference between (1.1) and $EX$ is bounded in absolute value by

$$(\frac{15}{2}e^{-m/b} + m^2Ce^{-\alpha b} + 4m^{\frac{1}{2}}e^{-m/2})\frac{m\log m}{E(K \wedge m)} + Ce^{-\alpha b}(\frac{m\log m}{E(K \wedge b)})^2.$$

*Remark.* If $P\{K > k\} \leq Ce^{-\alpha k}$ for all $k$, then letting $b \approx \sqrt{m}$ in the Proposition 4.5 bound shows that the approximation error converges to zero faster than $\exp(-m^{1/3})$ as $m \to \infty$.

**Proof.**

$$\max_{j \geq 1}\sum_{r=j}^{m-1}\frac{m}{m-r}P\{K > r | K \geq j\} \leq m\sum_{r=1}^{m-1}\frac{1}{m-r} < m\log m.$$

Also,

$$\sum_{r=b}^{m-1}\frac{m}{m-r}P\{K > r\}\sum_{j=1}^{r}\frac{m}{m-j+1} \leq Ce^{-\alpha b}(m\log m)^2,$$

16

$$\sum_{i=0}^{b-1} \frac{m}{m-i} P\{K > i\} \geq E(K \wedge b).$$

and

$$\sum_{i=0}^{m-1} \frac{m}{m-i} P\{K > i\} \geq E(K \wedge m).$$

Combining these bounds and Proposition 4.2 with Proposition 4.3 proves Proposition 4.5.

$\square$

*Remark.* If $m = 800$ and the $K_i$'s are geometric $(\frac{1}{2})$, then taking $b = 40$ in Proposition 4.5 yields the error bound 0.0016. The next proposition is much more effective for geometric $K_i$'s.

**Proposition 4.6.** If the hazard function $h_a$ of the $K_i$'s is bounded below by $\delta > 0$ for $a < m$, and if $(m - b)\delta \geq 1$, then

$$\|P\{A_\infty \in \cdot\} - P\{A_\infty^{(b)} \in \cdot\}\| \leq 6(m-b)^{\frac{1}{2}} \left\{ \frac{\delta^\delta}{(1+\delta)^{1+\delta}} \right\}^{m-b} + 8m^{\frac{1}{2}} e^{-m/2} + 2m^2(1-\delta)^b.$$

and the difference between (1.1) and $EX$ is bounded in absolute value by

$$\left[ 3(m-b)^{\frac{1}{2}} \left\{ \frac{\delta^\delta}{(1+\delta)^{1+\delta}} \right\}^{m-b} + 4m^{\frac{1}{2}} e^{-m/2} + m^2(1-\delta)^b \right] \frac{m(1-\delta)}{\delta} + (1-\delta)^b \left( \frac{m \log m}{E(K \wedge b)} \right)^2.$$

Proposition 4.6 is proved in Sellke (1992).

$\square$

*Remark.* If the $K_i$'s are geometric $(\frac{1}{2})$ and $m = 100$, then taking $b = 62$ and $\delta = \frac{1}{2}$ causes the bounds in Proposition 4.6 to be $1.08 \times 10^{-14}$ and $5.64 \times 10^{-13}$.

## 5. Generalization

Suppose we start with $r$ red balls and $w = m - r$ white balls in the box. Let $Y_{r,\ell}$ be the number of iid samples needed to see (or paint red) $\ell - r$ of the white balls, so that we go from $r$ red balls at the start to $\ell$ red balls at the end. (Thus, $X = Y_{0,m}$.) The formula

$$(5.1) \qquad \frac{\sum_{i=r}^{\ell-1} \frac{1}{m-i}}{\sum_{i=0}^{m-1} \frac{1}{m-i} P\{K > i\}} + \frac{\sum_{i=1}^{m-1} \frac{1}{m-i} P\{K > i\} \sum_{j=1}^{i} \frac{1}{m-j+1}}{\left[ \sum_{i=0}^{m-1} \frac{1}{m-i} P\{K > i\} \right]^2}$$

17

should (usually) be a good approximation for $EY_{r,\ell}$ when the first term of (5.1) is large. The argument for the first term is exactly the same as in Section 2: the numerator is the expected number of single-ball draws needed to get the required number of white balls, and the denominator is the expected number of draws needed to complete a sample. The second term of (5.1) (which is exactly the same as the second term of (1.1)) is again a correction for the error in the first term caused by "boundary overshoot." In this case, the $\{(A_n, U_n)\}_{n=1}^{\infty}$ chain of Section 3 starts either in state $(1, w)$ or in state $(1, w - 1)$, depending on whether the first ball chosen is red or white. The $(A_n, U_n)$ chain is absorbed by the states on $U$-level $m - \ell$. If a successful coupling can be achieved with high probability between this $(A_n, U_n)$ chain and an $(A_n^{(b)}, U_n^{(b)})$ chain, then the $A$ coordinate of the absorbing state will have a distribution approximated by (3.12). Providing also that the bound in Lemma 4.3 is small, the second term of (5.1) will do a good job of correcting the error in the first term.

If the sample sizes are geometric $(1/2)$ and the box contains $m = 30$ balls, $r = 18$ of which are red, then the expected value $EY_{18,26}$ of the number of samples need to see $26 - 18 = 8$ of the white balls is $15.289268$, while (5.1) gives $15.289282$, for an error of $0.000014$. For the same situation with $K_i$'s uniformly distributed on $\{1, 2, 3, 4, 5\}$, $EY_{18,26} = 10.18018794$, while (5.1) gives $10.18018721$, for an error of $-0.00000073$.

**References**

Lindvall, T. (1992), *Lectures on the Coupling Method*, Wiley, New York.

Feller, W. (1968), *An Introduction to Probability Theory and Its Applications*, Volume 1, 3rd ed., Wiley, New York.

Polya, G. (1930), Eine wahrscheinlichkeitsaufgabe zur kundenwerbung. *Z. Angew. Math. Mech.*, **10** (1930), 96-97. *In George Pólya: Collected Papers* (1984), Vol IV, MIT Press, Gian-Carlo Rota, Cambridge.

Ross, S. (1988), *A First Course in Probability*, 3rd ed., Macmillan, New York.

Sellke, T. (1992), How many iid samples does it take to see all the balls in a box?, Purdue University Techinical Report 92-47.